



Docket No.: PF-0525 USN

Certificate of Mailing

I hereby certify that the enclosed correspondence is being deposited with the United States Postal Service as first class mail in an envelope addressed to:  
Commissioner for Patents, P.O. Box 1450, Alexandria, VA 22313-1450 on 6/23/03

By: [Signature]

Printed: DELLS

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of: Lal et al.

Title: HUMAN SOCS PROTEINS

Serial No.: 09/701,232

Filing Date: July 5, 2001

Examiner: Hamud, F.

Group Art Unit: 1647

Commissioner for Patents  
P.O. Box 1450  
Alexandria, VA 22313-1450

**DECLARATION OF LARS MICHAEL FURNESS  
UNDER 37 C.F.R. § 1.132**

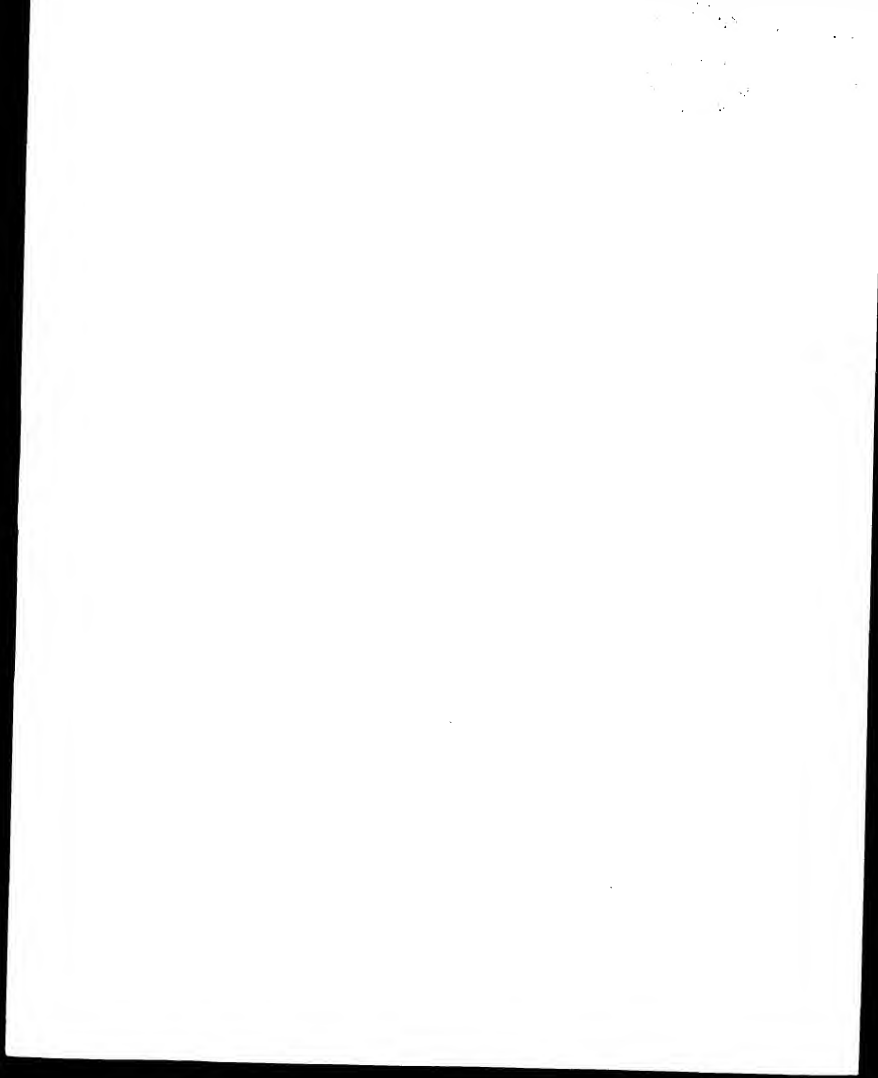
I, L. MICHAEL FURNESS, a citizen of the United Kingdom, residing at 2 Brookside, Exning, Newmarket, United Kingdom, declare that:

1. I was employed by Incyte Corporation (hereinafter "Incyte") as a Director of Pharmacogenomics until December 31, 2001. I am currently under contract to be a Consultant to Incyte Corporation.

2. In 1984, I received a B.Sc.(Hons) in Biomolecular Science (Biophysics and Biochemistry) from Portsmouth Polytechnic.

From 1985-1987 I was at the School of Pharmacy in London, United Kingdom, during which time I analyzed lipid methyltransferase enzymes using a variety of protein analysis methods, including one-dimensional (1D) and two-dimensional (2D) gel electrophoresis, HPLC, and a variety of enzymatic assay systems.

#11  
895  
7/16/03



I then worked in the Protein Structure group at the National Institute for Medical Research until 1989, setting up core facilities for nucleic acid synthesis and sequencing, as well as assisting in programs on protein kinase C inhibitors.

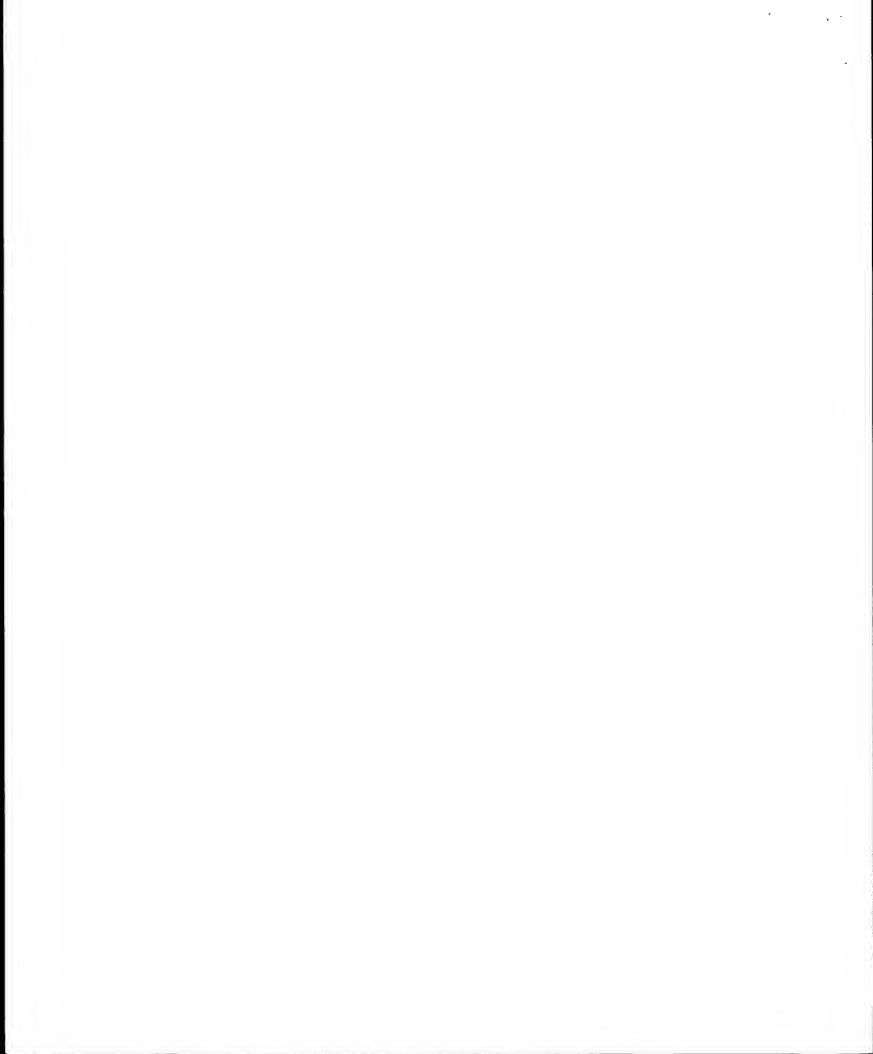
After a year at Perkin Elmer-Applied Biosystems as a technical specialist, I worked at the Imperial Cancer Research Fund between 1990-1992, on a Eureka-funded program collaborating with Amersham Pharmacia in the United Kingdom and CEPH (Centre d'Etude du Polymorphisme Humaine) in Paris, France, to develop novel nucleic acid purification and characterization methods.

In 1992, I moved to Pfizer Central Research in the United Kingdom, where I stayed until 1998, initially setting up core DNA sequencing and then a DNA arraying facility for gene expression analysis in 1993. My work also included bioinformatics, and I was responsible for the support of all Pfizer neuroscience programs in the United Kingdom. This then led me into carrying out detailed bioinformatics and wet lab work on the sodium channels, including antibody generation, western and northern analyses, PCR, tissue distribution studies, and sequence analyses on novel sequences identified.

In 1998, I moved to Incyte Genomics, Inc., to the Pharmacogenomics group to look at the application of genomics and proteomics to the pharmaceutical industry. In 1999, I was appointed director of the LifeExpress Lead Program which used microarray and protein expression data to identify pharmacologically and toxicologically relevant mechanisms to assist in improved drug design and development.

On December 12, 2001 I founded Nuomics Consulting Ltd., in Exning, U.K., and I am currently employed as Managing Director. Nuomics Consulting Ltd. provides expert technical knowledge and advice to businesses around the areas of genomics, proteomics, pharmacogenomics, toxicogenomics and chemogenomics.

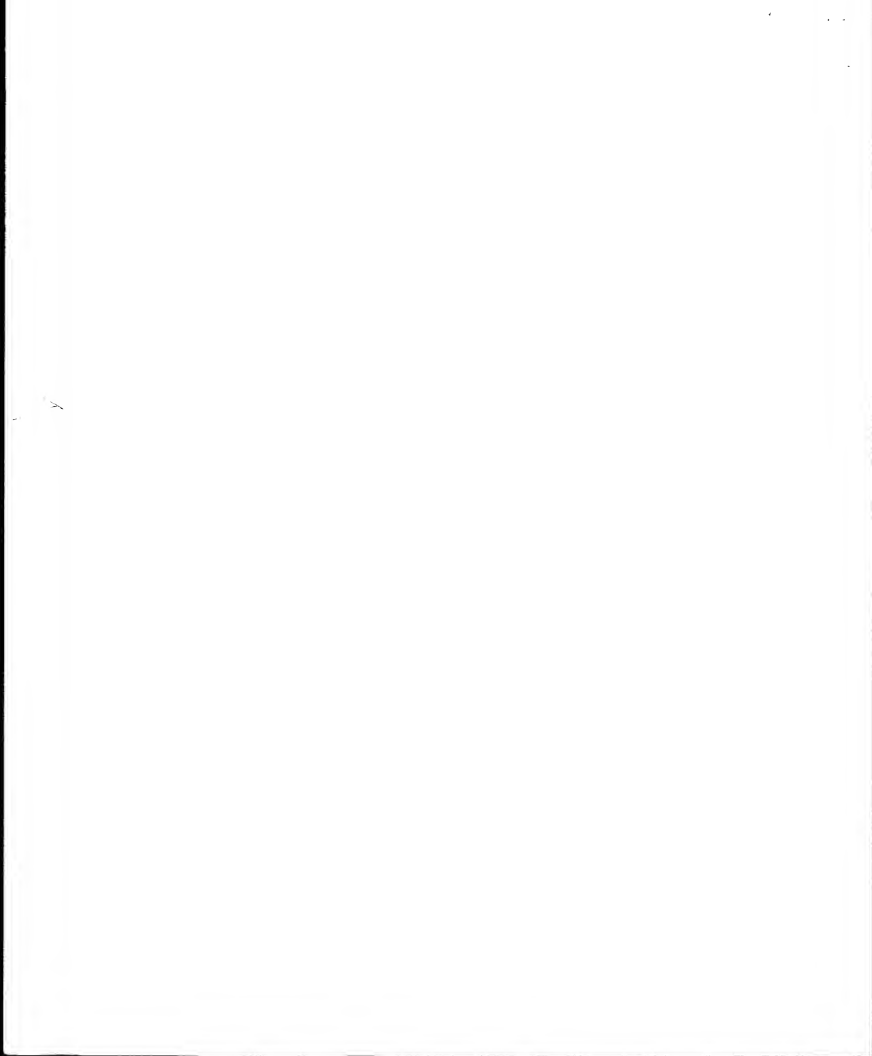
3. I have reviewed the specification of a United States patent application that I understand was filed on July 5, 2001 in the names of Preeti Lal et al. and was assigned Serial No. 09/701,232 (hereinafter "the Lal '232 application"). Furthermore, I understand that this United States patent application is the National Stage of International Application No. PCT/US99/11497, filed May 25, 1999, and published in English as WO 99/61614 on December



2, 1999, which claims the benefit under 35 U.S.C. § 119(e) of provisional applications U.S. Ser. No. 60/087,104, filed May 28, 1998 (hereinafter the Lal '104 application) and U.S. Ser. No. 60/150,701, filed December 17, 1998. The provisional applications provide support for what is disclosed in the instant Lal '232 application. The SEQ ID NO:5 and SEQ ID NO:14 sequences recited in the Lal '232 application claims were first disclosed in the Lal '104 application and listed as SEQ ID NO:5 and SEQ ID NO:11, respectively, in the Lal '104 application. My remarks herein will therefore be directed to the Lal '104 patent application, and May 28, 1998, as the relevant date of filing. In broad overview, the Lal '104 specification pertains to certain nucleotide and amino acid sequences and their use in a number of applications, including gene and protein expression monitoring applications that are useful in connection with (a) developing drugs (e.g., for the treatment of cancer), and (b) monitoring the activity of drugs for purposes relating to evaluating their efficacy and toxicity.

4. I understand that (a) the Lal '232 application contains claims that are directed to a isolated polypeptide comprising the amino acid sequence of SEQ ID NO:5 (hereinafter "the SEQ ID NO:5 polypeptide"), and (b) the Patent Examiner has rejected those claims on the grounds that the specification of the Lal '232 application does not disclose a specific and substantial asserted utility or a well established utility for the claimed SEQ ID NO:5 polypeptide. I further understand that whether or not a patent specification discloses a substantial, specific and credible utility for its claimed subject matter is properly determined from the perspective of a person skilled in the art to which the specification pertains at the time of the patent application was filed. In addition, I understand that a substantial, specific and credible utility under the patent laws must be a "real-world" utility.

5. I have been asked (a) to consider with a view to reaching a conclusion (or conclusions) as to whether or not I agree with the Patent Examiner's position that the Lal '232 application and its priority application, the Lal '104 application, do not disclose a specific and substantial asserted utility or a well established "real-world" utility for the claimed SEQ ID NO:5 polypeptide, and (b) to state and explain the bases for any conclusions I reach. I have been informed that, in connection with my considerations, I should determine whether or not a person

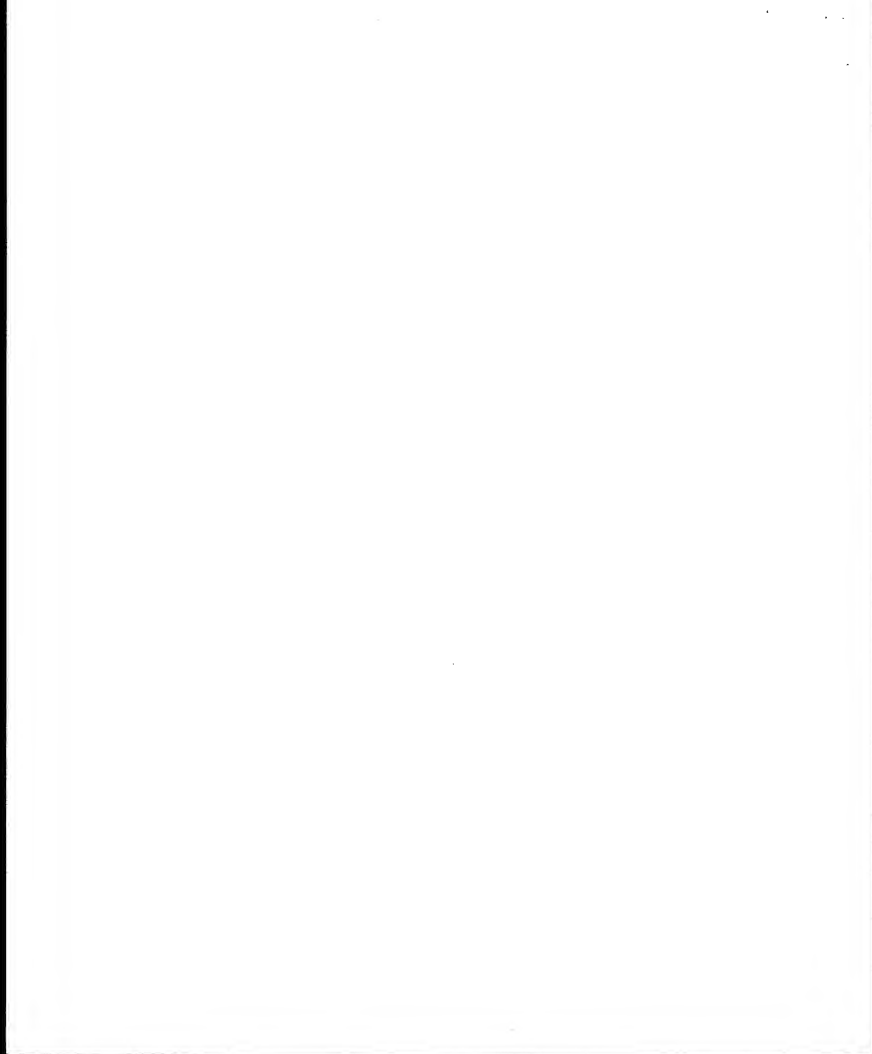


skilled in the art to which the Lal '104 application pertains on May 28, 1998, would have concluded that the Lal '104 application disclosed, for the benefit of the public, a specific beneficial use of the SEQ ID NO:5 polypeptide in its then available and disclosed form. I have also been informed that, with respect to the "real-world" utility requirement, the Patent and Trademark Office instructs its Patent Examiners in Section 2107.01 of the Manual of Patent Examining Procedure, 8<sup>th</sup> Edition, August 2001, under the heading I. Specific and Substantial Requirements, Research Tools):

Many research tools such as gas chromatographs, screening assays, and nucleotide sequencing techniques have a clear, specific and unquestionable utility (e.g., they are useful in analyzing compounds). An assessment that focuses on whether an invention is useful only in a research setting thus does not address whether the specific invention is in fact "useful" in a patent sense. Instead, Office personnel must distinguish between inventions that have a specifically identified substantial utility and inventions whose asserted utility requires further research to identify or reasonably confirm.

6. I have considered the matters set forth in paragraph 5 of this Declaration and have concluded that, contrary to the position I understand the Patent Examiner has taken, the specification of the Lal '104 patent application disclosed to a person skilled in the art at the time of its filing a number of substantial, specific and credible real-world utilities for the claimed SEQ ID NO:5 polypeptide. More specifically, persons skilled in the art on May 28, 1998 would have understood the Lal '104 application to disclose the use of the SEQ ID NO:5 polypeptide as a research tool in a number of gene and protein expression monitoring applications that were well-known at that time to be useful in connection with the development of drugs and the monitoring of the activity of such drugs. I explain the bases for reaching my conclusion in this regard in paragraphs 7-13 below.

7. In reaching the conclusion stated in paragraph 6 of this Declaration, I considered (a) the specification of the Lal '104 application, and (b) a number of published articles and patent documents that evidence gene and protein expression monitoring techniques that were well-known before the May 28, 1998 filing date of the Lal '104 application. The published articles and patent documents I considered are:





- (a) Anderson, N.L., Esquer-Blasco, R., Hofmann, J.-P., Anderson, N.G., A Two-Dimensional Gel Database of Rat Liver Proteins Useful in Gene Regulation and Drug Effects Studies, Electrophoresis, 12, 907-930 (1991) (hereinafter "the Anderson 1991 article") (copy annexed at Tab A);
- (b) Anderson, N.L., Esquer-Blasco, R., Hofmann, J.-P., Mehues, L., Raymackers, J., Steiner, S. Witzmann, F., Anderson, N.G., An Updated Two-Dimensional Gel Database of Rat Liver Proteins Useful in Gene Regulation and Drug Effect Studies, Electrophoresis, 16, 1977-1981 (1995) (hereinafter "the Anderson 1995 article") (copy annexed at Tab B);
- (c) Wilkins, M.R., Sanchez, J.-C., Gooley, A.A., Appel, R.D., Humphery-Smith, I., Hochstrasser, D.F., Williams, K.L., Progress with Proteome Projects: Why all Proteins Expressed by a Genome Should be Identified and How To Do It, Biotechnology and Genetic Engineering Reviews, 13, 19-50 (1995) (hereinafter "the Wilkins article") (copy annexed at Tab C);
- (d) Celis, J.E., Rasmussen, H.H., Leffers, H., Madsen, P., Honore, B., Gesser, B., Dejgaard, K., Vandekerckhove, J., Human Cellular Protein Patterns and their Link to Genome DNA Sequence Data: Usefulness of Two-Dimensional Gel Electrophoresis and Microsequencing, FASEB Journal, 5, 2200-2208 (1991) (hereinafter "the Celis article") (copy annexed at Tab D);
- (e) Franzen, B., Linder, S., Okuzawa, K., Kato, H., Auer, G., Nonenzymatic Extraction of Cells from Clinical Tumor Material for Analysis of Gene Expression by Two-Dimensional Polyacrylamide Gel Electrophoresis, Electrophoresis, 14, 1045-1053 (1993) (hereinafter "the Franzen article") (copy annexed at Tab E);
- (f) Bjellqvist, B., Basse, B., Olsen, E., Celis, J.E., Reference Points for Comparisons of Two-Dimensional Maps of Proteins from Different Human Cell Types Defined in a pH Scale Where Isoelectric Points Correlate with Polypeptide Compositions, Electrophoresis, 15, 529-539 (1994) (hereinafter "the Bjellqvist article") (copy annexed at Tab F); and
- (g) Large Scale Biology Company Info; LSB and LSP Information; from <http://www.lsb.com> (2001) (copy annexed at Tab G).

the 1990s, the number of people in the UK who are aged 65 and over has increased from 10.5 million to 12.5 million, and the number of people aged 75 and over has increased from 4.5 million to 6.5 million (Office of National Statistics 2000).

There is a growing awareness of the need to address the needs of older people in the community. The Department of Health (1999) has published a strategy for older people, which sets out a vision for the future of older people's health and social care. The strategy is based on the following principles: older people should be able to live independently and actively in the community; older people should be able to access the services they need; and older people should be able to participate in decisions about their care and services.

The strategy also sets out a number of key objectives for the future of older people's health and social care. These include: to improve the health and well-being of older people; to ensure that older people have access to the services they need; to ensure that older people are able to participate in decisions about their care and services; and to ensure that older people are able to live independently and actively in the community.

The strategy is a key document for the future of older people's health and social care in the UK. It sets out a vision for the future of older people's health and social care, and sets out a number of key objectives for the future of older people's health and social care. The strategy is a key document for the future of older people's health and social care in the UK.

The strategy is a key document for the future of older people's health and social care in the UK. It sets out a vision for the future of older people's health and social care, and sets out a number of key objectives for the future of older people's health and social care. The strategy is a key document for the future of older people's health and social care in the UK.

The strategy is a key document for the future of older people's health and social care in the UK. It sets out a vision for the future of older people's health and social care, and sets out a number of key objectives for the future of older people's health and social care. The strategy is a key document for the future of older people's health and social care in the UK.

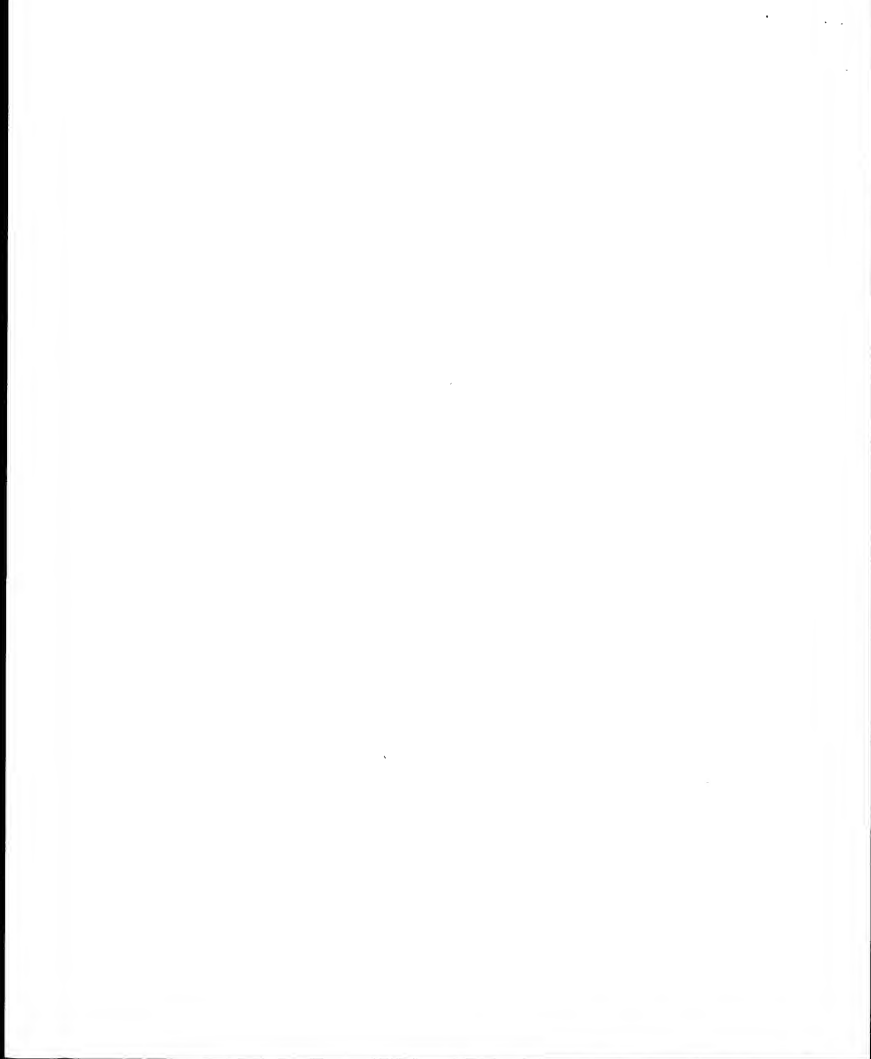
The strategy is a key document for the future of older people's health and social care in the UK. It sets out a vision for the future of older people's health and social care, and sets out a number of key objectives for the future of older people's health and social care. The strategy is a key document for the future of older people's health and social care in the UK.

The strategy is a key document for the future of older people's health and social care in the UK. It sets out a vision for the future of older people's health and social care, and sets out a number of key objectives for the future of older people's health and social care. The strategy is a key document for the future of older people's health and social care in the UK.

8. Many of the published articles I considered (i.e., at least items (a)-(f) identified in paragraph 7) relate to the development of protein two-dimensional gel electrophoretic techniques for use in protein expression monitoring applications in drug development and toxicology. As I will discuss below, a person skilled in the art who read the Lal '104 application on May 28, 1998 would have understood that application to disclose the SEQ ID NO:5 polypeptide to be useful for a number of gene and protein expression monitoring applications, e.g., in the use of two-dimensional polyacrylamide gel electrophoresis and western blot analysis of tissue samples in drug development and in toxicity testing.

Furthermore, items (a)-(f) establish that protein two-dimensional polyacrylamide gel electrophoresis and western blot analysis were well-known and established methods routinely used in toxicology testing and drug development at the time of filing the Lal '104 application and for several years prior to May 28, 1998. As such, one of ordinary skill in the art would have recognized that the polypeptide of SEQ ID NO:5 could be used in toxicology testing and drug development, irrespective of its biochemical activities.

9. The SEQ ID NO:5 and SEQ ID NO:14 sequences recited in the Lal '232 application claims were first disclosed in the Lal '104 application and listed as SEQ ID NO:5 and SEQ ID NO:11, respectively, in the Lal '104 application. The SEQ ID NO:5 polypeptide is referred to as HSCOP-5 in the Lal '232 application and as SOCP-5 in the Lal '104 application. Turning more specifically to the Lal '104 specification, the SEQ ID NO:5 polypeptide is shown at pages 46-47 under the heading "Sequence Listing." The Lal '104 specification specifically teaches that the "invention features substantially purified polypeptides, human SOCS proteins, referred to collectively as 'SOCP' and individually as 'SOCP-1', 'SOCP-2', 'SOCP-3', 'SOCP-4', 'SOCP-5', and 'SOCP-6' and that the "invention provides a substantially purified polypeptide comprising an amino acid sequence selected from the group consisting of SEQ ID NO:1, SEQ ID NO:2, SEQ ID NO:3, SEQ ID NO:4, SEQ ID NO:5, SEQ ID NO:6 (SEQ ID NO:1 through 6). . . " (Lal '104 application at page 2, lines 32-36). It further teaches that (a) the identity of the SEQ ID NO:5 polypeptide was determined from a uterus tissue cDNA library (UTRSNOR01) (Lal '104 application, Tables 1 and 4), (b) the SEQ ID NO:5 polypeptide is the human SOCS protein referred to as "SOCP-5" and is encoded by SEQ ID NO:11. (Lal '104 application at page 2, lines

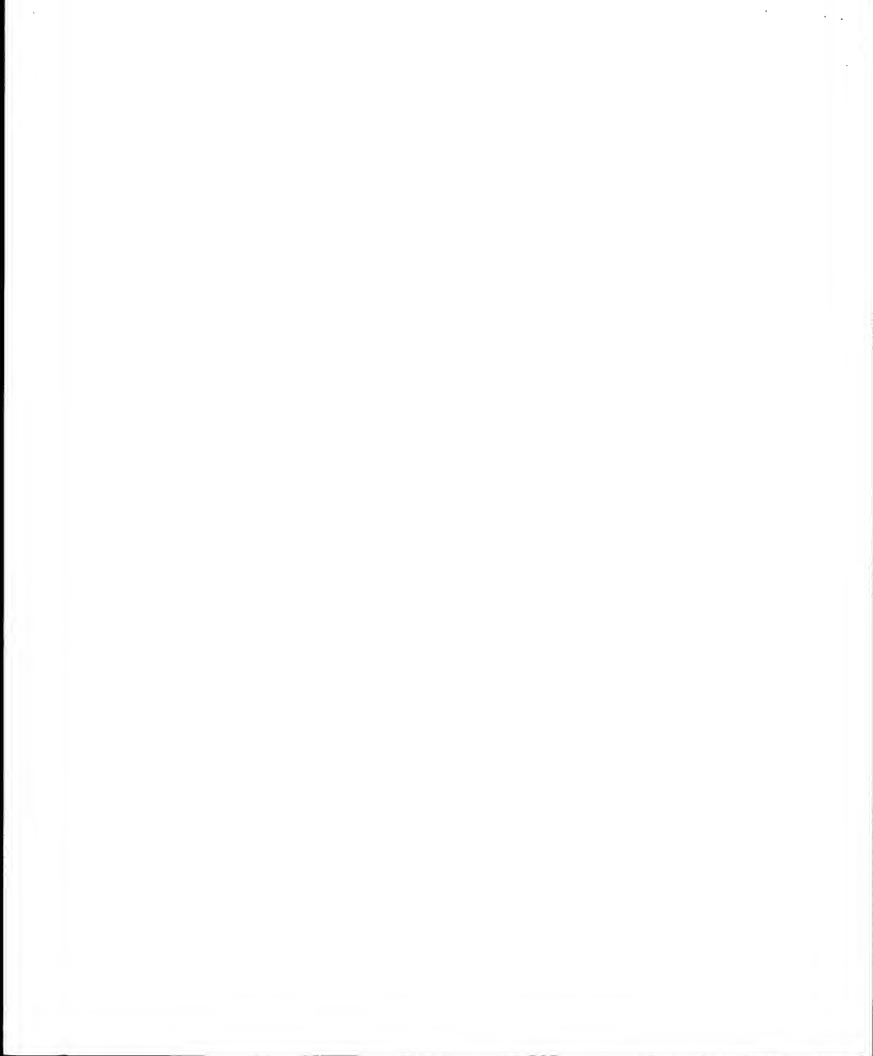


31-36 and Table 1), and (c) northern analysis of SEQ ID NO:11 shows its expression predominantly in cDNA libraries made from reproductive, cardiovascular, hematopoietic/immune, cancer-associated, inflammation-associated, and fetal tissues (Lal '104 application at Table 3) and therefore "SOCP appears to play a role in cancer, immune disorders, and infectious diseases." (Lal '104 application at page 20, lines 22-23.)

The Lal '104 application discusses a number of uses of the SEQ ID NO:5 polypeptide in addition to its use in gene and protein expression monitoring applications. I have not fully evaluated these additional uses in connection with the preparation of this Declaration and do not express any views in this Declaration regarding whether or not the Lal '104 specification discloses these additional uses to be substantial, specific and credible real-world utilities of the SEQ ID NO:5 polypeptide. Consequently, my discussion in this Declaration concerning the Lal '104 application focuses on the portions of the application that relate to the use of the SEQ ID NO:5 polypeptide in gene and protein expression monitoring applications.

10. The Lal '104 application discloses that the polynucleotide sequences disclosed therein, including the polynucleotides encoding the SEQ ID NO:5 polypeptide, are useful as probes in chip based technologies. It further teaches that the chip based technologies can be used "for the detection and/or quantification of nucleic acid or protein sequences." (Lal '104 application at page 18, lines 27-28.)

The Lal '104 application also discloses that the SEQ ID NO:5 polypeptide is useful in other protein expression detection technologies. The Lal '104 application states that "[i]mmunological methods for detecting and measuring the expression of SOCP using either specific polyclonal or monoclonal antibodies are known in the art. Examples of such techniques include enzyme-linked immunosorbent assays (ELISAs), radioimmunoassays (RIAs), and fluorescence activated cell sorting (FACS)." (Lal '104 application at page 18, lines 29-32.) Furthermore, the Lal '104 application discloses that "[a] variety of protocols for measuring SOCP, including ELISAs, RIAs, and FACS, are known in the art and provide a basis for diagnosing altered or abnormal levels of SOCP expression. Normal or standard values for SOCP expression are established by combining body fluids or cell extracts taken from normal

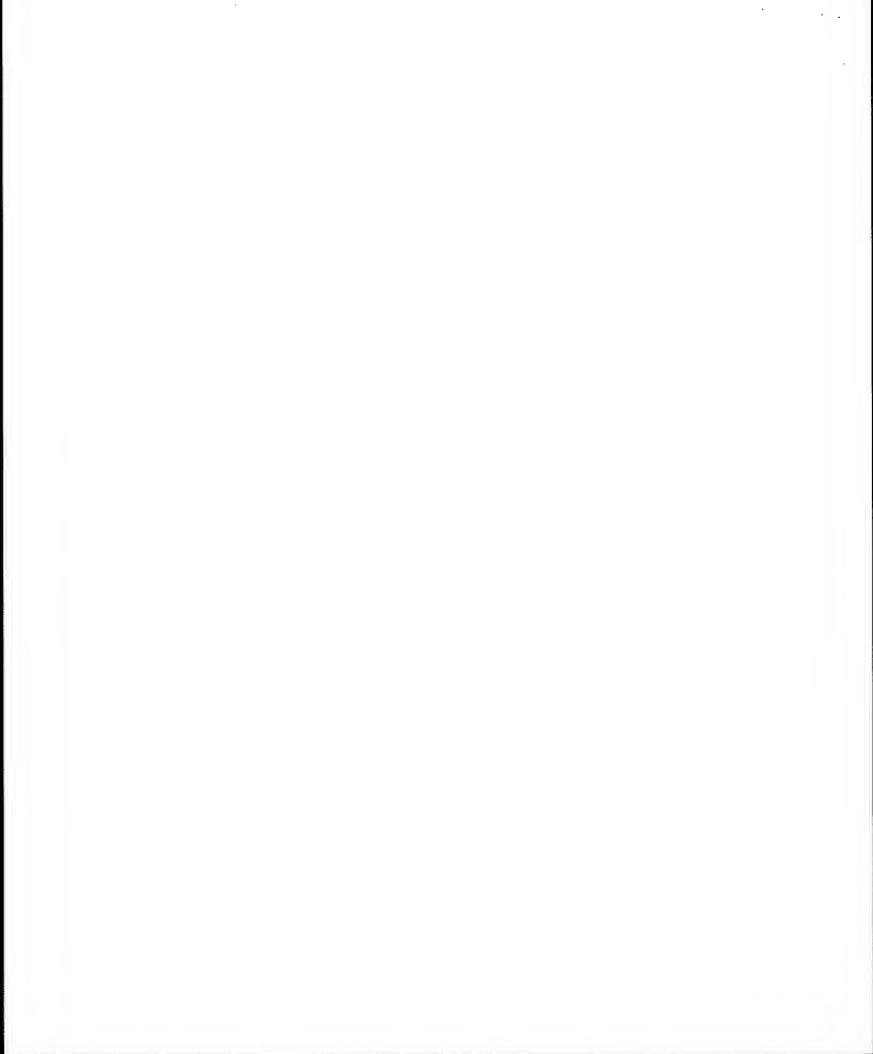


mammalian subjects, preferably human, with antibody to SOCP under conditions suitable for complex formation." (Lal '104 application at page 28, lines 5-9.)

In addition, at the time of filing the Lal '104 application, it was well known in the art that gene and protein expression analyses also included two-dimensional polyacrylamide gel electrophoresis (2-D PAGE) technologies, which were developed during the 1980s, and as exemplified by the Anderson 1991 and 1995 articles (Tab A and Tab B). The Anderson 1991 article teaches that a 2-D PAGE map has been used to connect and compare hundreds of 2-D gels of rat liver samples from a variety of studies including regulation of protein expression by various drugs and toxic agents (Tab A at page 907). The Anderson 1991 article teaches an empirically-determined standard curve fitted to a series of identified proteins based upon amino acid chain length (Tab A at page 911) and how that standard curve can be used in protein expression analysis. The Anderson 1991 article teaches that "there is a long-term need for a comprehensive database of liver proteins" (Tab A at page 912).

The Wilkins article is one of a number of documents that were published prior to the May 28, 1998 filing date of the Lal '104 application that describes the use of the 2-D PAGE technology in a wide range of gene and protein expression monitoring applications, including monitoring and analyzing protein expression patterns in human cancer, human serum plasma proteins, and in rodent liver following exposure to toxins. In view of the Lal '104 application, the Wilkins article, and other related pre-May 28, 1998 publications, persons skilled in the art on May 28, 1998 clearly would have understood the Lal '104 application to disclose the SEQ ID NO:5 polypeptide to be useful in 2-D PAGE analyses for the development of new drugs and monitoring the activities of drugs for such purposes as evaluating their efficacy and toxicity, as explained more fully in paragraph 12 below.

With specific reference to toxicity evaluations, those of skill in the art who were working on drug development on May 28, 1998 (and for many years prior to May 28, 1998) without any doubt appreciated that the toxicity (or lack of toxicity) of any proposed drug they were working on was one of the most important criteria to be considered and evaluated in connection with the development of the drug. They would have understood at that time that good drugs are not only potent, they are specific. This means that they have strong effects on a specific biological target and minimal effects on all other biological targets. Ascertaining that a





candidate drug affects its intended target, and identification of undesirable secondary effects (i.e., toxic side effects), had been for many years among the main challenges in developing new drugs. The ability to determine which genes are positively affected by a given drug, coupled with the ability to quickly and at the earliest time possible in the drug development process identify drugs that are likely to be toxic because of their undesirable secondary effects, have enormous value in improving the efficiency of the drug discovery process, and are an important and essential part of the development of any new drug. In fact, the desire to identify and understand toxicological effects using the experimental assays described above led Dr Leigh Anderson to found the Large Scale Biology Corporation in 1985, in order to pursue commercial development of the 2-D electrophoretic protein mapping technology he had developed. In addition, the company focused on toxicological effects on the proteome as clearly demonstrated by its goals and by its senior management credentials described in company documents (see Tab G at pages 1, 3, and 5).

Accordingly, the teachings in the Lal '104 application, in particular regarding use of SEQ ID NO:5 in differential gene and protein expression analysis (2-D PAGE maps) and in the development and the monitoring of the activities of drugs, clearly include toxicity studies and persons skilled in the art who read the Lal '104 application on May 28, 1998 would have understood that to be so.

11. As previously discussed (*supra*, paragraphs 7 and 8), my experience with protein analysis methods in the mid-1980s and the several publications annexed to this Declaration at Tabs A through F evidence information that was available to the public regarding two-dimensional polyacrylamide gel electrophoresis technology and its uses in drug discovery and toxicology testing before the May 28, 1998 filing date of the Lal '104 application. In particular the Celis article stated that "protein databases are expected to foster a variety of biological information.... -- among others, ..... drug development and testing" (See Tab D, page 2200, second column). The Franzen article shows that 2-D PAGE maps were used to identify proteins in clinical tumor material (See Tab E). The Lal '104 application clearly discloses that expression of SOCP-5 is associated with reproductive, cardiovascular, hematopoietic/immune, cancer-associated, inflammation-associated, and fetal tissues (Lal '104 application at Table 3). The Bjellqvist article showed that a protein may be identified accurately by its positional co-



ordinates, namely molecular mass and isoelectric point (See Tab F). The Lal '104 application clearly disclosed SEQ ID NO:5 from which it would have been routine for one of skill in the art to predict both the molecular mass and the isoelectric point using algorithms well known in the art at the time of filing.

12. A person skilled in the art on May 28, 1998, who read the Lal '104 application, would understand that application to disclose the SEQ ID NO:5 polypeptide to be highly useful in analysis of differential expression of proteins. For example, the specification of the Lal '104 application would have led a person skilled in the art on May 28, 1998 who was using protein expression monitoring in connection with working on developing new drugs for the treatment of cancer, immune disorders, and infectious diseases to conclude that a 2-D PAGE map that used the isolated SEQ ID NO:5 polypeptide would be a highly useful tool and to request specifically that any 2-D PAGE map that was being used for such purposes utilize the SEQ ID NO:5 polypeptide sequence. Expressed proteins are useful for 2-D PAGE analysis in toxicology expression studies for a variety of reasons, particularly for purposes relating to providing controls for the 2-D PAGE analysis, and for identifying sequence or post-translational variants of the expressed sequences in response to exogenous compounds. Persons skilled in the art would appreciate that a 2-D PAGE map that utilized the SEQ ID NO:5 polypeptide sequence would be a more useful tool than a 2-D PAGE map that did not utilize this protein sequence in connection with conducting protein expression monitoring studies on proposed (or actual) drugs for treating cancer, immune disorders, and infectious diseases for such purposes as evaluating their efficacy and toxicity.

I discuss in more detail in items (a)-(b) below a number of reasons why a person skilled in the art, who read the Lal '104 specification on May 28, 1998, would have concluded based on that specification and the state of the art at that time, that SEQ ID NO:5 polypeptide would be a highly useful tool for analysis of a 2-D PAGE map for evaluating the efficacy and toxicity of proposed drugs for cancer, immune disorders, and infectious diseases by means of 2-D PAGE maps, as well as for other evaluations:

(a) The Lal '104 specification contains a number of teachings that would lead persons skilled in the art on May 28, 1998 to conclude that a 2-D PAGE map that utilized

the 1990s, the number of people in the UK who are employed in the public sector has increased by 1.5 million, from 2.5 million in 1980 to 4 million in 1995. The public sector has become an important employer of people with mental health problems.

There is a growing awareness of the need to improve the mental health of people in the public sector. The Department of Health (1996) has published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'. The strategy also states that 'the Department of Health will work with other government departments to ensure that the mental health of people in the public sector is given the same priority as the physical health of people in the public sector'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

The Department of Health has also published a strategy for mental health care, which includes a commitment to improve the mental health of people in the public sector. The strategy states that 'the mental health of people in the public sector is a priority for the Department of Health'.

the isolated SEQ ID NO:5 polypeptide would be a more useful tool for protein expression monitoring applications relating to drugs for treating cancer, immune disorders, and infectious diseases than a 2-D PAGE map that did not use the SEQ ID NO:5 polypeptide sequence. Among other things, the Lal '104 specification teaches that (i) the identity of the SEQ ID NO:5 polypeptide was determined from a "uterus tissue cDNA library (UTRSNOR01)," (Lal '104 application, Tables 1 and 4) (ii) the SEQ ID NO:5 polypeptide is the human SOCS protein referred to as "SOCP-5" (listed as HSCOP-5 in the Lal '232 application) (Lal '104 application at page 2, lines 31-36 and Table 1), and (iii) SEQ ID NO:11 (listed as SEQ ID NO:14 in the Lal '232 application) is expressed predominantly in cDNA libraries made from reproductive, cardiovascular, hematopoietic/immune, cancer-associated, inflammation-associated, and fetal tissues (Lal '104 application at Table 3) and therefore "SOCP appears to play a role in cancer, immune disorders, and infectious diseases." (Lal '104 application at page 20, lines 22-23; see paragraph 9, *supra*). The isolated polypeptide could therefore be used as a control to more accurately gauge the expression of SOCP-5 (listed as HSCOP-5 in the Lal '232 application) in the sample and consequently more accurately gauge the affect of a toxicant on expression of the gene.

(b) Persons skilled in the art on May 28, 1998 would have appreciated (i) that the protein expression monitoring results obtained using a 2-D PAGE map that utilized a SEQ ID NO:5 polypeptide would vary, depending on the particular drug being evaluated, and (ii) that such varying results would occur both with respect to the results obtained from the SEQ ID NO:5 polypeptide and from the 2-D PAGE map as a whole (including all its other individual proteins). These kinds of varying results, depending on the identity of the drug being tested, in no way detracts from my conclusion that persons skilled in the art on May 28, 1998, having read the Lal '104 specification, would specifically request that any 2-D PAGE map that was being used for conducting protein expression monitoring studies on drugs for treating cancer, immune disorders, and infectious diseases (*e.g.*, a toxicology study or any efficacy study of the type that typically takes place in connection with the development of a drug) utilize the SEQ ID NO:5 polypeptide sequence. Persons skilled in the art on May 28, 1998 would have wanted their 2-D PAGE map to utilize the SEQ ID NO:5 polypeptide sequence because a 2-D PAGE map that utilized protein sequence information the polypeptide (as compared to one that did not) would



provide more useful results in the kind of protein expression monitoring studies using 2-D PAGE maps that persons skilled in the art have been doing since well prior to May 28, 1998.

The foregoing is not intended to be an all-inclusive explanation of all my reasons for reaching the conclusions stated in this paragraph 12, and in paragraph 6, *supra*. In my view, however, it provides more than sufficient reasons to justify my conclusions stated in paragraph 6 of this Declaration regarding the Lal '104 application disclosing to persons skilled in the art at the time of its filing substantial, specific and credible real-world utilities for the SEQ ID NO:5 polypeptide.

13. Also pertinent to my considerations underlying this Declaration is the fact that the Lal '104 disclosure regarding the uses of the SEQ ID NO:5 polypeptide for protein expression monitoring applications is not limited to the use of that protein in 2-D PAGE maps. For one thing, the Lal '104 disclosure regarding the technique used in gene and protein expression monitoring applications is broad. (Lal '104 application at, e.g., page 18, lines 24-28.)

In addition, the Lal '104 specification repeatedly teaches that the protein described therein (including the SEQ ID NO:5 polypeptide) may desirably be used in any of a number of long established "standard" techniques, such as ELISA or western blot analysis, for conducting protein expression monitoring studies. See, e.g.:

(a) Lal '104 application at page 18, lines 29-32 ("Immunological methods for detecting and measuring the expression of SOCP using either specific polyclonal or monoclonal antibodies are known in the art. Examples of such techniques include enzyme-linked immunosorbent assays (ELISAs), radioimmunoassays (RIAs), and fluorescence activated cell sorting (FACS).");

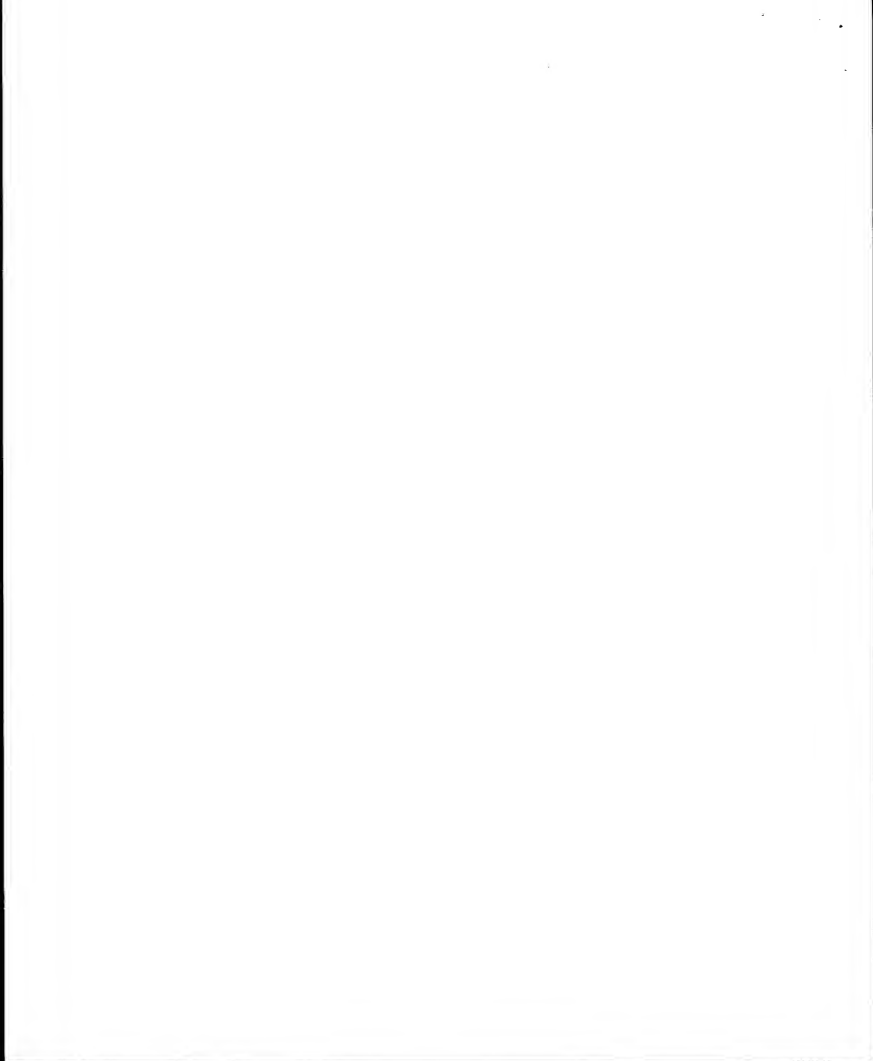
(b) Lal '104 application at page 28, lines 5-12 ("A variety of protocols for measuring SOCP, including ELISAs, RIAs, and FACS, are known in the art and provide a basis for diagnosing altered or abnormal levels of SOCP expression. Normal or standard values for SOCP expression are established by combining body fluids or cell extracts taken from normal mammalian subjects, preferably human, with antibody to SOCP under conditions suitable for complex formation. The amount of standard complex formation may be quantitated by various methods, preferably by photometric means. Quantities of SOCP expressed in subject, control,





and disease samples from biopsied tissues are compared with the standard values. Deviation between standard and subject values establishes the parameters for diagnosing disease.”).

Thus a person skilled in the art on May 28, 1998, who read the Lal ‘104 specification, would have routinely and readily appreciated that the SEQ ID NO:5 polypeptide disclosed therein would be useful to conduct protein expression monitoring analyses using 2-D PAGE mapping or western blot analysis or any of the other traditional membrane-based protein expression monitoring techniques that were known and in common use many years prior to the filing of the Lal ‘104 application. For example, a person skilled in the art on May 28, 1998 would have routinely and readily appreciated that the SEQ ID NO:5 polypeptide would be a useful tool in conducting protein expression analyses, using the 2-D PAGE mapping or western analysis techniques, in furtherance of (a) the development of drugs for the treatment of cancer, immune disorders, and infectious diseases, and (b) analyses of the efficacy and toxicity of such drugs.



14. I declare further that all statements made herein of my own knowledge are true and that all statements made herein on information and belief are believed to be true; and further, that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, and that willful false statements may jeopardize the validity of this application and any patent issuing thereon.

---

L. Michael Furness, B.Sc.

Signed at Exning, United Kingdom

this \_\_\_\_ day of \_\_\_\_\_, 2003

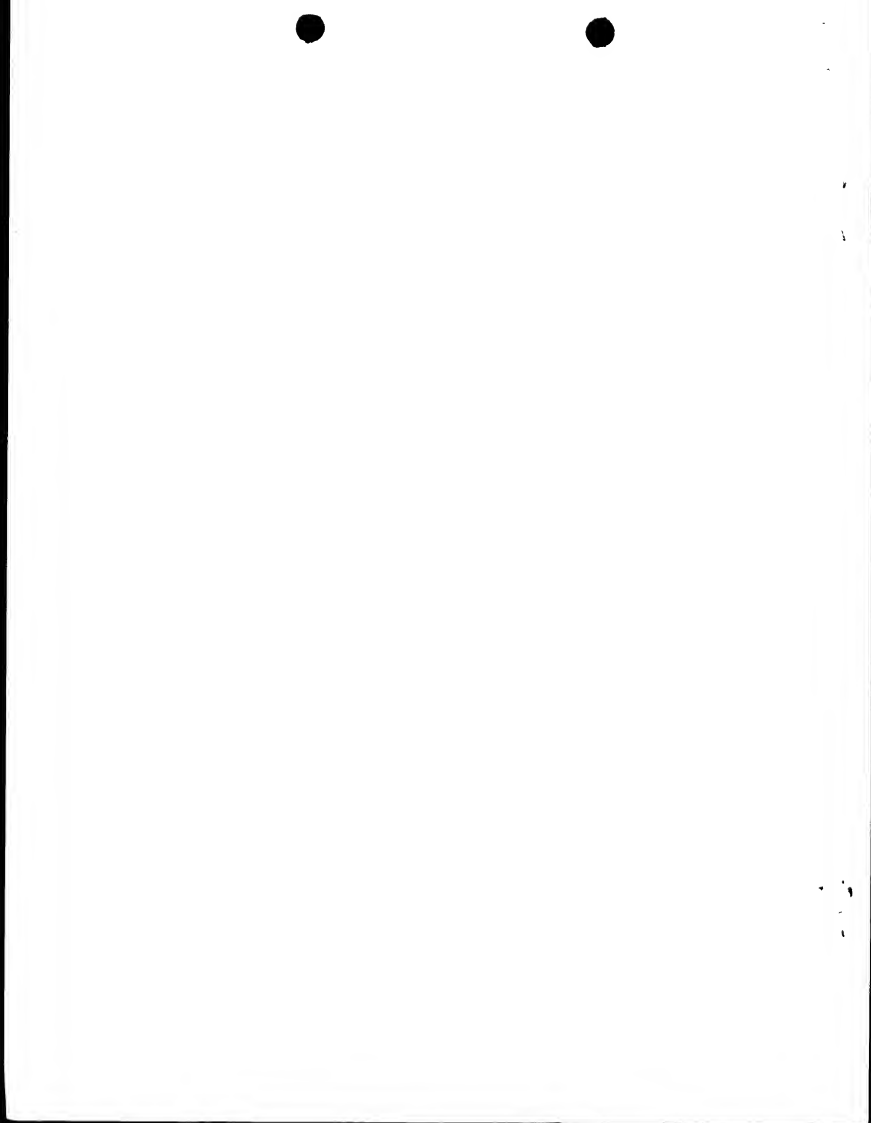




wkt

## INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification <sup>6</sup> : C12N 15/12, C07K 14/47, A61K 38/17, G01N 33/68, C12Q 1/68, C07K 16/18		A3	(11) International Publication Number: <b>WO 99/61614</b>
(21) International Application Number: PCT/US99/11497		(43) International Publication Date: 2 December 1999 (02.12.99)	
(22) International Filing Date: 25 May 1999 (25.05.99)			
(30) Priority Data: 60/087,104 28 May 1998 (28.05.98) US 09/216,006 17 December 1998 (17.12.98) US		(74) Agents: BILLINGS, Lucy, J. et al.; Incyte Pharmaceuticals, Inc., 3174 Porter Drive, Palo Alto, CA 94304 (US).	
(63) Related by Continuation (CON) or Continuation-in-Part (CIP) to Earlier Applications US 60/087,104 (CIP) Filed on 28 May 1998 (28.05.98) US 09/216,006 (CIP) Filed on 17 December 1998 (17.12.98)		(81) Designated States: AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, GM, HR, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZW, ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).	
(71) Applicant (for all designated States except US): INCYTE PHARMACEUTICALS, INC. [US/US]; 3174 Porter Drive, Palo Alto, CA 94304 (US).		Published <i>With international search report.</i> <i>Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.</i>	
(72) Inventors; and (75) Inventors/Applicants (for US only): LAL, Preeti [IN/US]; 2382 Lass Drive, Santa Clara, CA 95054 (US); HILLMAN, Jennifer, L. [US/US]; 230 Monroe Drive #12, Mountain View, CA 94040 (US); GORGONE, Gina [US/US]; 1253 Pinecrest Drive, Boulder Creek, CA 95006 (US); CORLEY, Neil, C. [US/US]; 1240 Dale Avenue #30, Mountain View, CA 94040 (US); PATTERSON, Chandra [US/US]; 490 Sherwood Way #1, Menlo Park, CA 94025 (US); YUE, Henry [US/US]; 826 Lois Avenue, Sunnyvale, CA 94087		(88) Date of publication of the international search report: 2 March 2000 (02.03.00)	
(54) Title: HUMAN SOCS PROTEIN'S			
(57) Abstract  The invention provides human SOCS proteins (HSCOP) and polynucleotides which identify and encode HSCOP. The invention also provides expression vectors, host cells, antibodies, agonists, and antagonists. The invention also provides methods for diagnosing, treating, or preventing disorders associated with expression of HSCOP.			



# FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	VU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						





# INTERNATIONAL SEARCH REPORT

International Application No.  
PCT/US 99/11497

A. CLASSIFICATION OF SUBJECT MATTER  
IPC 6 C12N15/12 C07K14/47 A61K38/17 G01N33/68 C12Q1/68  
C07K16/18

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)  
IPC 6 C07K C12N A61K G01N C12Q

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>WO 98 20023 A (INST MEDICAL W &amp; E HALL ;VINEY ELIZABETH M (AU); STARR ROBYN (AU);) 14 May 1998 (1998-05-14) see SEQ ID NO: 24-27 (pp. 142-147) see the claims abstract; examples 5-8,11,18-24,28; table 7.1 page 4 -page 5 page 17 -page 18 page 33</p> <p style="text-align: center;">---</p> <p style="text-align: center;">-/--</p>	1-16,19

☒ Further documents are listed in the continuation of box C.

☒ Patent family members are listed in annex.

### \* Special categories of cited documents:

- \*A\* document defining the general state of the art which is not considered to be of particular relevance
- \*E\* earlier document but published on or after the international filing date
- \*L\* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- \*O\* document referring to an oral disclosure, use, exhibition or other means
- \*P\* document published prior to the international filing date but later than the priority date claimed

\*T\* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

\*X\* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

\*Y\* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

\*S\* document member of the same patent family

Date of the actual completion of the international search

21 September 1999

Date of mailing of the international search report

13. 01. 00

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentaan 2  
NL - 2280 HV Rijswijk  
Tel (+31-70) 340-2040, Tx. 31 651 epo nl,  
Fax (+31-70) 340-3016

Authorized officer

Oderwald, H



# INTERNATIONAL SEARCH REPORT

International Application No.

PCT/US 99/11497

## C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	<p>DATABASE EMEST16 [Online]  EMBL, Heidelberg, Germany  AC: AA401503, ID: H51200297,  29 April 1997 (1997-04-29)  HILLIER L ET AL.: "Homo sapiens cDNA clone  742641"  XP002115960  abstract</p> <p style="text-align: center;">---</p>	3-13
X	<p>WO 92 19734 A (INDIANA UNIVERSITY  FOUNDATION ;UNIV YALE (US))  12 November 1992 (1992-11-12)  see SEQ ID NO: 33 and 34 (pp.145-151)  abstract; claims  1,21,31,33,63-65,75,84,95,99,103,111,119;  figure 24  page 17 -page 19</p> <p style="text-align: center;">---</p>	3-14,16
A	<p>D J HILTON ET AL: "Twenty proteins  containing a C-terminal SOCS box form five  structural classes"  PROCEEDINGS OF THE NATIONAL ACADEMY OF  SCIENCES OF USA,  vol. 95, 1 January 1998 (1998-01-01),  pages 114-119, XP002085497  ISSN: 0027-8424  cited in the application  the whole document</p> <p style="text-align: center;">-----</p>	



## INTERNATIONAL SEARCH REPORT

International application No.

PCT/US 99/11497

**Box I Observations where certain claims were found unsearchable (Continuation of item 1 of first sheet)**

This International Search Report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1. ☒ Claims Nos.:  
because they relate to subject matter not required to be searched by this Authority, namely:

See FURTHER INFORMATION sheet PCT/ISA/210

2. ☒ Claims Nos.:  
because they relate to parts of the International Application that do not comply with the prescribed requirements to such an extent that no meaningful International Search can be carried out, specifically:

See FURTHER INFORMATION sheet PCT/ISA/210

3. ☐ Claims Nos.:  
because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

**Box II Observations where unity of invention is lacking (Continuation of item 2 of first sheet)**

This International Searching Authority found multiple inventions in this international application, as follows:

See additional sheet.

1. ☐ As all required additional search fees were timely paid by the applicant, this International Search Report covers all searchable claims.
2. ☐ As all searchable claims could be searched without effort justifying an additional fee, this Authority did not invite payment of any additional fee.
3. ☐ As only some of the required additional search fees were timely paid by the applicant, this International Search Report covers only those claims for which fees were paid, specifically claims Nos.:
4. ☒ No required additional search fees were timely paid by the applicant. Consequently, this International Search Report is restricted to the invention first mentioned in the claims, it is covered by claims Nos.:

Claims 1-16, 19 (all partially)

Remark on Protest

☐ The additional search fees were accompanied by the applicant's protest.

☐ No protest accompanied the payment of additional search fees.



**FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210**

Continuation of Box 1.2

Claims Nos.: 17, 18, 20

Claims 17, 18, 20 have not been searched due to insufficient disclosure of the claimed compounds.

The applicant's attention is drawn to the fact that claims, or parts of claims, relating to inventions in respect of which no international search report has been established need not be the subject of an international preliminary examination (Rule 66.1(e) PCT). The applicant is advised that the EPO policy when acting as an International Preliminary Examining Authority is normally not to carry out a preliminary examination on matter which has not been searched. This is the case irrespective of whether or not the claims are amended following receipt of the search report or during any Chapter II procedure.





FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210

1. Claims: 1-16, 19 (all partially)

A substantially purified polypeptide comprising amino acid sequence SEQ ID NO. 1 and fragments thereof, a variant having at least 90% identity; an isolated and purified polynucleotide encoding said polypeptide; a variant of said polynucleotide having at least 90% identity; a polynucleotide which hybridizes under stringent conditions to said polynucleotide; a polynucleotide having a sequence which is complementary to said polynucleotide; a method for detecting a polynucleotide encoding said polypeptide; said method wherein the polynucleotide is amplified by applying PCR; an isolated and purified polynucleotide comprising polynucleotide sequence SEQ ID NO. 10 and fragments thereof, or a variant having at least 90% identity; a polynucleotide having a sequence which is complementary to said polynucleotide; an expression vector comprising at least a fragment of said polynucleotide; a host cell comprising said expression vector; a method for producing a polypeptide comprising amino acid sequence SEQ ID NO. 1; a pharmaceutical composition comprising said polypeptide in conjunction with a suitable pharmaceutical carrier; an antibody which specifically binds to said polypeptide.

2. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 2 and 11.

3. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 3 and 12.

4. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 4 and 13.

5. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 5 and 14.

6. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 6 and 15.

7. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 7 and 16.



**FURTHER INFORMATION CONTINUED FROM PCT/ISA/ 210**

8. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 8 and 17.

9. Claims: 1-16, 19 (all partially)

Same as subject 1 but limited to SEQ ID NOS. 9 and 18.



# INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No

PCT/US 99/11497

Patent document cited in search report		Publication date	Patent family member(s)	Publication date
WO 9820023	A	14-05-1998	AU 4694397 A	29-05-1998
			EP 0948522 A	13-10-1999
			GB 2331753 A	02-06-1999
			NO 992116 A	29-06-1999
-----				
WO 9219734	A	12-11-1992	AU 675203 B	30-01-1997
			AU 1919792 A	21-12-1992
			CA 2102208 A	12-11-1992
			EP 0576623 A	05-01-1994
			EP 0933082 A	04-08-1999
			EP 0930365 A	21-07-1999
			EP 0930366 A	21-07-1999
			JP 7503123 T	06-04-1995
			US 5648464 A	15-07-1997
			US 5849869 A	15-12-1998
			US 5856441 A	05-01-1999
			US 5789195 A	04-08-1998
-----				



**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						





N. Leigh Anderson  
Ricardo Esquer-Blasco  
Jean-Paul Hofmann  
Norman G. Andersson

Large Scale Biology Corporation,  
Rockville, MD

## A two-dimensional gel database of rat liver proteins useful in gene regulation and drug effects studies

A standard two-dimensional (2-D) protein map of Fischer 344 rat liver (F344MST3) is presented, with a tabular listing of more than 1200 protein species. Sodium dodecyl sulfate (SDS) molecular mass and isoelectric point have been established, based on positions of numerous internal standards. This map has been used to connect and compare hundreds of 2-D gels of rat liver samples from a variety of studies, and forms the nucleus of an expanding database describing rat liver proteins and their regulation by various drugs and toxic agents. An example of such a study, involving regulation of cholesterol synthesis by cholesterol-lowering drugs and a high-cholesterol diet, is presented. Since the map has been obtained with a widely used and highly reproducible 2-D gel system (the Iso-Dalt<sup>®</sup> system), it can be directly related to an expanding body of work in other laboratories.

### Contents

1 Introduction	907
2 Material and methods	908
2.1 Sample preparation	908
2.2 Two-dimensional electrophoresis	909
2.3 Staining	909
2.4 Positional standardization	909
2.5 Computer analysis	910
2.6 Graphical data output	910
2.7 Experiment LSBC04	910
3 Results and discussion	910
3.1 The rat liver protein 2-D map	910
3.2 Carbamylated charge standards computed pI's and molecular mass standardization	911
3.3 An example of rat liver gene regulation: Cholesterol metabolism	911
3.3.1 MSN 413 (putative cytosolic HMG-CoA synthase) and sets of spots regulated coordinately or inversely	911
3.3.2 MSN 235 and correlated spots	912
3.3.3 An example of an anti-synergistic effect	912
3.3.4 Complexity of the cholesterol synthesis pathway	912
4 Conclusions	912
5 References	914
6 Addendum 1: Figures 1-13	914
Addendum 2: Tables 1-4	923
Table 1. Master table of proteins in rat liver database	923
Table 2. Table of some identified proteins	928
Table 3. Computed pI's of two sets of carbamylated protein standards: rabbit muscle CPK and human Hb.	929
Table 4. Computed pI's of some known proteins related to measured CPK pI's	930

### 1 Introduction

High-resolution two-dimensional electrophoresis of proteins, introduced in 1975 by O'Farrell and others [1-4], has been used over the ensuing 16 years to examine a wide variety of biological systems, the results appearing in more than 5000 published papers. With the advent of computerized systems for analyzing two-dimensional (2-D) gel images and constructing spot databases, it is also possible to plan and assemble integrated bodies of information describing the appearance and regulation of thousands of protein gene products [5, 6]. Creating such databases involves amassing and organizing quantitative data from thousands of 2-D gels, and requires a substantial commitment in technology and resources.

Given the long-term effort required to develop a protein database, the choice of a biological system takes on considerable importance. While *in vitro* systems are ideal for answering many experimental questions, especially in cancer research and genetics, our experience with cell cultures and tissue samples suggests that some *in vivo* approaches could have major advantages. In particular, we have noticed that liver tissue samples from rats and mice appear to show greater quantitative reproducibility (in terms of individual protein expression) than replicate cell cultures. This is perhaps a natural result of the homeostasis maintained in a complete animal vs. the well-known variability of cell cultures, the latter due principally to differences in reagents (e.g., fetal bovine serum), conditions (e.g., pH) and genetic "evolution" of cell lines while in culture. It is also more difficult to generate adequate amounts of protein from cell culture systems (particularly with attached cells), forcing the investigator to resort to radioisotope-based or silver-based detection methods. While these methods are more sensitive (sometimes much more sensitive) than the Coomassie Brilliant Blue (CBB) stain typically used for protein detection in "large" protein samples, they are generally more variable, more labor-intensive and, in the case of radiographic methods, may generate highly "noisy" images, due to the properties of the films used. By contrast, large protein samples can easily be prepared from liver using urea/Nonidet P-40 (NP-40) solubilization and stained with CBB, which has the advantage of being easily reproducible [8]. Finally, there remains the question of the "truthfulness" of many *in vitro* systems as compared to their *in vivo* analogs; how great are the changes caused by the introduction into a cul-

Correspondence: Dr. N. Leigh Anderson, Large Scale Biology Corporation, 9620 Medical Center Drive, Rockville, MD 20850, USA

Abbreviations: CBB, Coomassie Brilliant Blue; CPK, creatine phosphokinase; 2-D, two-dimensional; IEF, isoelectric focusing; MSN, master spot number; NP-40, Nonidet P-40; SDS, sodium dodecyl sulfate

© 1991 Verlagsgesellschaft mbH, D-6940 Weinheim, 1991

0173-0835/91/1111-0907\$3.50+2.50

ture and the associated shift to strong selection for growth, and how do these affect experimental outcomes? Hence the apparent advantages of *in vitro* systems, in terms of experimental manipulation, may be counterbalanced by other factors relating to 2-D data quality.

There is a second important class of reasons for exploring the use of an *in vivo* biological system such as the liver. Historically, there have been two broad approaches to the mechanistic dissection of biochemical processes in intact cellular systems: genetics (a search for informative mutants) and the use of chemical agents (drugs and chemical toxins). Both approaches help us to understand complex systems by disrupting some specific functional element and showing us the result. With the development of techniques for genetic manipulation and cloning, the genetic approach can be effectively applied either *in vitro* or *in vivo*, although the *in vitro* route is usually quicker. The chemical approach can also be applied to either sort of biological system; here, however, the bulk of consistently acquired information is in experimental animals (rats and mice). While most biologists know a short list of compounds having specific, experimentally useful effects (e.g., inhibitors of protein synthesis, ionophores, polymerase inhibitors, channel blockers, nucleotide analogs, and compounds affecting polymerization of cytoskeletal proteins), there is a much larger number of interesting chemically-induced effects, most of them characterized by toxicologists and pharmacologists in rodent systems. Just as a thorough genetic analysis would involve saturating a genome with mutations, it is possible to imagine a saturating number of drugs, the analysis of whose actions would reveal the complete biochemistry of the cell. While organized drug discovery efforts usually target specific desired effects, the nature of the process, with its dependence on screening large numbers of compounds, necessarily produces many unanticipated effects. It is therefore reasonable to suppose that the required broad range of compounds necessary to achieve "biochemical saturation" may be forthcoming; in fact, it may already exist among the hundreds of thousands of compounds that failed to qualify as drugs.

Among organs, the liver is an obvious choice for the study of chemical effects because of its well-known plasticity and responsiveness. The brain appears to be quite plastic (e.g. [7]), but it is a complicated mixture of cell types requiring skillful dissection for most experiments. The kidney, while quite responsive, also presents a potentially confounding mixture of cell types. The liver, by contrast, is made up of one predominant cell type which is easy to solubilize: the hepatocyte, representing more than 95% of its mass. Most importantly, the liver performs many homeostatic functions that require rapid modulation of gene expression. It appears that most chemical agents tested affect gene expression in the liver at some dosage (N. Leigh Anderson, unpublished observations), an interesting contrast to our earlier work with lymphocytes, for example, which seem to be much less responsive. Such results conform to the expectation that cells with a homeostatic, physiological role should be more plastic than cells differentiated for a purpose dependent on the action of a limited number of specific genes.

The liver also allows the parallels between *in vitro* and *in vivo* systems to be examined in detail. Significant progress

has been made in the development of mouse, rat and human hepatocyte culture systems, as well as in precision-cut tissue slices. Using such an array of techniques, it is possible to assemble a matrix of mammalian systems including mouse and rat *in vivo* on one level and mouse, rat and human *in vitro* on a second level, and to compare effects between species and between systems. This approach allows us to draw informed conclusions regarding the biochemical "universality" of biological responses among the mammalian and to offer some insight into the validity of *in vitro* approaches for toxicological screening. We believe this data will be necessary if *in vitro* alternatives are to achieve wide usage in government-mandated safety testing of drugs, consumer products and industrial and agricultural chemicals.

A number of interesting studies have been published using 2-D mapping to examine effects in the rodent liver. A number of investigators have made use of the technique to screen for existing genetic variants [8-11] or induced mutations [12-14], mainly in the mouse. This work builds on the wealth of genetic information available on the mouse and its established position as a mammalian mutation-detection system. While some studies of chemical effects have been undertaken in the mouse [15-17], most have used the rat [18-23]. The examination of the cytochrome p-450 system, in particular, has been carried out almost exclusively on the rat [24, 25].

These considerations lead us to conclude that rodent liver offers the best opportunity to systematically examine an array of gene regulation systems, and ultimately to build a predictive model of large-scale mammalian gene control. The basic underlying foundation of such a project is a reliable, reproducible master 2-D pattern of liver, to which ongoing experimental results can be referred. In this paper, we report such a master pattern for the acidic and neutral proteins of rat liver (pattern F344MST3). In future, this master will be supplemented by maps of basic proteins, and analogous maps of mouse and human liver.

## 2 Materials and methods

### 2.1 Sample preparation

Liver is an ideal sample material for most biochemical studies, including 2-D analysis. A sample is taken of approximately 0.5 g of tissue from the apical end of the left lobe of the liver. Solubilization is effected as rapidly as practical; a delay of 5-15 min appears to cause no major alteration in liver protein composition if the liver pieces are kept cold (e.g., on ice) in the interim. In the solubilization process, the liver sample is weighed, placed in a glass homogenizer (e.g., 15 mL Wheaton); 8 volumes of solubilizing solution\*

\* The solubilizing solution is composed of 2% NP-40 (Sigma), 9% urea (analytical grade, e.g. BDH or Bio-Rad), 0.5% dithiothreitol (DTT; Sigma) and 2% carrier ampholytes (pH 9-11 LKB; these come as a 20% stock solution, so 2% final concentration is achieved by making the final solution 10% 9-11 Ampholine by volume). A large batch of solubilant (several hundred mL) is made and stored frozen at -80°C in aliquots sufficient to provide enough for one day's estimated sample preparation requirement. The solution is never allowed to become warmer than room temperature at any stage during preparation or thawing for use, since heating of concentrated urea solutions can produce contaminants that covalently modify proteins producing artefactual charge shifts. Once thawed, any unused solubilizer is discarded.

ed (i.e., 4 mL per 0.5 g tissue) and the mixture is homogenized using first the loose- and then the tight-fit glass pestle. This takes approximately 5 strokes with the pestle and is carried out at room temperature because it would crystallize out in the cold. Once the liver sample is thoroughly homogenized in the solubilizer, it is assumed that all the proteins are denatured (by the chaotropic effect of the urea and NP-40 detergent) and the enzymes inactivated by the high pH (~9.5). Therefore these samples may be kept at room temperature until they can be centrifuged frozen as a group (within several hours of preparation). The samples are centrifuged for  $6 \times 10^4$  g min (e.g., 500 000 g for 12 min using a Beckman TL-100 centrifuge). The centrifuge rotor is maintained at just below room temperature (e.g., 15–20°C), but not too cold, so as to prevent the precipitation of urea. The centrifuge of choice is a Beckman L-100 because of the sample tube sizes available, but any ultracentrifuge accepting smallish tubes will suffice. When an appropriate centrifuge is not available near the site of sample preparation, samples can be frozen at -80°C and thawed prior to centrifugation and collection of supernatants. Each supernatant is carefully removed following centrifugation and aliquoted into at least 4 clean tubes for storage. This is done by transferring all the supernatant to one clean tube, mixing this gently (to assure homogeneous composition) and then dividing it into 4 aliquots. The aliquots are frozen immediately at -80°C. These multiple aliquots can provide insurance against a failed run or a freezer meltdown.

## 2. Two-dimensional electrophoresis

Sample proteins are resolved by 2-D electrophoresis using the  $20 \times 25$  cm Iso-Dalt® 2-D gel system [26–29], produced by LSB and by Hoefer Scientific Instruments, San Francisco) operating with 20 gels per batch. All first-dimensional isoelectric focusing (IEF) gels are prepared using the same single standardized batch of carrier ampholytes BDH 4–8A in the present case, selected by LSB's batch-testing program for rat and mouse database work<sup>\*\*\*</sup>. A 10  $\mu$ L sample of solubilized liver protein is applied to each gel, and the gels are run for 33 000 to 34 500 volt-hours using a progressively increasing voltage protocol implemented by a programmable high-voltage power supply. An "Angelique" computer-controlled gradient-casting system (produced by LSB) is used to prepare second-dimensional sodium dodecyl sulfate (SDS) polyacrylamide gradient slab gels in which the top 5% of the gel is 11%T acrylamide, and the lower 95% of the gel varies linearly from 11% to 18%T.

This system has recently been modified so as to employ a commercially available 30.8%T acrylamide/*N,N*-methylenebisacrylamide prepared solution (thus avoiding the handling of the solid acrylamide monomer) and three additional stock solutions: buffer (made from Sigma pre-set Tris), persulfate and *N,N,N,N*-tetramethylethylenediamine (TEMED). Each gel is identified by a computer-printed filter paper label polymerized into the lower left corner of the gel. First-dimensional IEF tube gels are loaded

directly (as extruded) onto the slab gels without equilibration, and held in place by polyester fabric wedges ("Wedges", produced by LSB) to avoid the use of hot agarose. Second-dimensional slab gels are run overnight, in groups of 20, in cooled DALT tanks (10°C) with buffer circulation. All run parameters, reagent source and lot information, and notations of deviation from expected results are entered by the technician responsible on a detailed, multi-page record of the experiment.

## 2.3 Staining

Following SDS-electrophoresis, slab gels are stained for protein using a colloidal Coomassie Blue G-250 procedure in covered plastic boxes, with 10 gels (totalling approximately 1 L of gel) per box. This procedure (based on the work of Neuhoff [30, 31]) involves fixation in 1.5 L of 50% ethanol and 2% phosphoric acid for 2 h, three 30 min washes, each in 2 L of cold tap water, and transfer to 1.5 L of 34% methanol, 17% ammonium sulfate and 2% phosphoric acid for 1 h, followed by the addition of a gram of powdered Coomassie Blue G-250 stain. Staining requires approximately 4 days to reach equilibrium intensities, whereupon gels are transferred to cool tap water and their surfaces rinsed to remove any particulate stain prior to scanning. Gels may be kept for several months in water with added sodium azide. The water washes remove ethanol that would dissolve the stain and render the system noncolloidal, with high backgrounds. The concentrated ammonium sulfate and methanol solution is diluted by equilibration with the water volume of the gels to automatically achieve the correct final concentrations for colloidal staining. Practical advantages of this staining approach can be summarized as follows: (i) the low, flat background makes computer evaluation of small spots (max OD < 0.02) possible, especially when using laser densitometry; (ii) up to 1500 spots can be reliably detected on many gels (e.g., rat liver) at loadings low enough to preserve excellent resolution; and (iii) reproducibility appears to be very good: at least several hundred spots have coefficients of reproducibility less than 15%. This value is at least as good as previous CBB methods, and significantly better than many silver stain systems.

## 2.4 Positional standardization

The carbamylated rabbit muscle creatine phosphokinase (CPK) standards [32] are purchased from Pharmacia and BDH. Amino acid compositions, and numbers of residues present in proteins used for internal standardization, are taken from the Protein Identification Resource (PIR) sequence database [33].

## 2.5 Computer analysis

Stained slab gels are digitized in red light at 134 micron resolution, using either a Molecular Dynamics laser scanner (with pixel sampling) or an Eikonix 78/99 CCD scanner. Raw digitized gel images are archived on high-density DAT tape (or equivalent storage media) and a greyscale video-print prepared from the raw digital image as hard-copy backup of the gel image. Gels are processed using the Kepler<sup>®</sup> software system (produced by LSB), a commercially available workstation-based software package built on

<sup>\*\*\*</sup> This material (succeeding certified batches of which are available from Hoefer Scientific Instruments) has the most linear pH gradient produced by any ampholyte tested except for the Pharmacia wide range which has an undesirable tendency to bind high-molecular weight acidic proteins, causing them to streak.

some of the principles of the earlier TYCHO system [34-41]. Procedure PROC008 is used to yield a spotlist giving position, shape and density information for each detected spot. This procedure makes use of digital filtering, mathematical morphology techniques and digital masking to remove the background, and uses full 2-D least-squares optimization to refine the parameters of a 2-D Gaussian shape for each spot. Processing parameters and file locations are stored in a relational database, while various log files detailing operation of the automatic analysis software are archived with the reduced data. The computed resolution and level of Gaussian convergence of each gel are inspected and archived for quality control purposes.

Experiment packages are constructed using the Kepler experimental definition database to assemble groups of 2-D patterns corresponding to the experimental groups (e.g., treated and control animals). Each 2-D pattern is matched to the appropriate "master" 2-D pattern (pattern F344MST3 in the case of Fischer 344 rat liver), thereby providing linkage to the existing rodent protein 2-D databases. The software allows experiments containing hundreds of gels to be constructed and analyzed as a unit, with up to 100 gels displayed on the screen at one time for comparative purposes and multiple pages to accommodate experiments of > 1000 gels. For each treatment, proteins showing significant quantitative differences vs. appropriate controls are selected using group-wise statistical parameters (e.g., Student's *t*-test, Kepler<sup>®</sup> procedure STUDENT). Proteins satisfying various quantitative criteria (such as  $P < 0.001$  difference from appropriate controls) are represented as highlighted spots onscreen or on computer-plotted protein maps and stored as spot populations (i.e., logical vectors) in a liver protein database. Quantitative data (spot parameters, statistical or other computed values) are stored as real-valued vectors in the database. Analysis of coregulation is performed using a Pearson product-moment correlation (Kepler procedure CORREL) to determine whether groups of proteins are coordinately regulated by any of the treatments. Such groups can be presented graphically on a protein map, and reported together with the statistical criteria used to assess the level of coregulation. Multivariate statistical analysis (e.g., principal components' analysis) is performed on data exported to SAS (SAS Institute).

## 2.6 Graphical data output

Graphical results are prepared in GKS and translated within Kepler<sup>®</sup> into output for any of a variety of devices. Linedrawing output is typically prepared as Postscript and printed on an Apple Laserwriter. Detailed maps presented here have been generated using an ultra-high-resolution Postscript-compatible Linotronic output device. Greyscale graphics are reproduced from the workstation screen using a Seikosha videoprinter. Patterns are shown in the standard orientation, with high molecular mass at the top and acidic proteins to the left.

## 2.7 Experiment LSB04

In the study described here 12-week-old Charles River male F344 rats were used. Diets were prepared at LSB, based on a Purina 5755M Basal Purified Diet. Lovastatin and cholestyramine were obtained as prescription pharma-

ceuticals, ground and mixed with the diet at concentrations of 0.075% and 1%, respectively. The high cholesterol diet was Purina 5801M-A (5% cholesterol plus 1% sodium cholate in the control diet). Animal work was carried out by Microbiological Associates (Bethesda, MD). Animals were acclimatized for one week on the control diet, fed test or control diets for one week, and sacrificed on day 8. Average daily doses of lovastatin and cholestyramine in appropriate groups were 37 mg/kg/day and 5 g/kg/day, respectively, based on the weight of the food consumed. Liver samples were collected and prepared for 2-D electrophoresis according to the standard liver protocol (homogenization in 8 volumes of 9 M urea, 2% NP-40, 0.5% dithiothreitol, 2% LKB pH 9-11 carrier ampholytes, followed by centrifugation for 30 min at 80000 × *g*). Kidney, brain and plasma samples were frozen. Gels were run as described above, and the data was analyzed using the Kepler<sup>®</sup> system. Gels were scaled, to remove the effect of differences in protein loading, by setting the summed abundances of a large number of matched spots equal for each gel (linear scaling).

## 3 Results and discussion

### 3.1 The rat liver protein 2-D map

F344MST3 is a standard 2-D pattern of rat liver proteins based on the Fischer 344 strain. This pattern was initiated from a single 2-D gel and extensively edited in an experiment comparing it to a range of protein loads, so as to include both small spots and well-resolved representations of high-abundance spots. More than 700 rat liver 2-D patterns have been matched to F344MST3 in a series of drug effects and protein characterization experiments, and numerous new spots (induced by specific drugs, for instance) have been added as a result. A modified version including additional spots present in the Sprague-Dawley outbred rat has also been developed (data not shown). Figure 1 shows a greyscale representation and Fig. 2 a schematic plot of the master pattern. More than 1200 spots are included, most of which are visible on typical gels loaded with 10 µL of solubilized liver protein prepared by the standard method and stained with colloidal Coomassie Blue. Master spot numbers (MSN's) have been assigned to all proteins, and appear in the following figures, each showing one quadrant of the pattern. Figure 3 shows the upper left (acidic, high molecular mass) quadrant, Fig. 4 the upper right (basic, high molecular mass) quadrant, Fig. 5 the lower left (acidic, low molecular mass) quadrant, and Fig. 6 the lower right (basic, low molecular mass) quadrant. The quadrants overlap as an aid to moving between them. The gel position (in 100 micron units), isoelectric point (relative to the CPK calibration curve in Fig. 8) are listed for each spot (Table 1). Because of the precision of the CPK-*pI* values, these parameters can be used to relate spot locations between gel systems more reliably than using *pI* measurements expressed as pH. A major objective of current studies is the identification of all major spots corresponding to known liver proteins, as well as rigorous definitions of subcellular organelle contents. Of particular interest to us is the parallel development of identifications in the rat and mouse liver maps, allowing detailed comparisons of gene expression effects in the two systems. The results of these studies will be presented systematically in a later edition of this database.

We include here a useful series of 22 orienting identifications as an aid to other users of the rat liver pattern (Table 1).

## 2. Carbamylated charge standards, computed pI's and molecular mass standardization

We have previously shown that the use of a system of close-spaced internal pI markers (made by carbamylating a basic protein) offers an accurate and workable solution to the problem of assigning positions in the pI dimension [32]. The same system, based on 36 protein species made by carbamylating rabbit muscle CPK, has been used here to assign pI's to rabbit muscle CPK and neutral proteins. The standards were electrophoresed with total liver proteins, and the standard spots added to a special version of the master pattern F344MST3. The gel X-coordinates of all liver protein spots lying within the CPK charge train were then transformed into CPK pI positions by interpolation between the positions of immediately adjacent standards (Table 1) using a Kepler<sup>†</sup> vector procedure.

It has proven possible to compute fairly accurate pI values for many proteins from the amino acid composition [42]. We have attempted here to test a further elaboration of this approach, in which we computed pI's for the CPK standards themselves, based on our knowledge of the rabbit muscle CPK sequence and the fact that adjacent members of the charge train typically differ by blockage of one additional lysine residue (Table 3). We compared these values to similar computed pI's for an additional set of carbamylated standards made from human hemoglobin beta chains and a series of rat liver and human plasma proteins of known position and sequence (Fig. 7, Table 4). The result demonstrates good concordance between these systems. Two proteins show significant deviations: liver fatty-acid binding protein (FABP; #1 in Table 4) and protein disulphide isomerase (#20 in the table). The FABP spot present on F344MST3 may represent a charge-modified version of a more basic parent spot closer to the expected pI, not resolved in the IEF/SDS gel. Of particular importance is the fact that, by comparing computed pI's of sequenced but unlocated gel proteins with the CPK pI's, we can assign a probable gel location without making any assumptions regarding the actual gel pH gradient. This offers a useful shortcut, given the vagaries of pH measurement on small diameter IEF gels. We have used this approach to compute the CPK pI's of all rat and mouse proteins in the PIR sequence database, as an aid to protein identification (data not shown).

In order to standardize SDS molecular weight (SDS-MW), we have used a standard curve fitted to a series of identified proteins (Fig. 8). Rather than using molecular mass *per se*, we have elected to use the number of amino acids in the polypeptide chain, as perhaps a better indication of the length of the SDS-coated rod that is sieved by the second dimension slab. The resulting values were multiplied by 100 (the weighted average mass of amino acids in sequenced proteins) to give predicted molecular masses. Because we use gradient slabs, we have not constrained the fit curve to conform to any predetermined model; rather we tried many equations and selected the best using the program "Tablecurve" on a PC. The equation chosen was  $y = a + bx + cx^2$ , where  $y$  is the number of residues,  $x$  is the gel

Y coordinate,  $a$  is 511.83,  $b$  is -0.2731 and  $c$  is 33183801. The resulting fit appears to be fairly good over a broad range of molecular mass.

## 3.3 An example of rat liver gene regulation: Cholesterol metabolism

Experiment LSBC04 was designed as a small-scale test of the regulation of cholesterol metabolism *in vivo* by three agents included in the diet: lovastatin (Mevacor<sup>®</sup>, an inhibitor of HMG-CoA reductase), cholestyramine (a bile acid sequestrant that has the effect of removing cholesterol from the gut-liver recirculation), and cholesterol itself. The first two agents should lower available cholesterol and the third should raise it, allowing manipulation of relevant gene expression control systems in both directions. Such an experiment offers an interesting test of the 2-D mapping system since most of the pathway enzymes are present in low abundance, many are membrane-bound and difficult to solubilize, and the pathway itself is complex. Approximately 1000 proteins were separated and detected in liver homogenates. Twenty-one proteins were found to be affected by at least one treatment, and these could be divided into several coregulated groups.

### 3.3.1 MSN 413 (putative cytosolic HMG-CoA synthase) and sets of spots regulated coordinately or inversely

One group of spots (including a spot assigned to the cytosolic HMG-CoA synthase, MSN 413) showed the expected increase in abundance with lovastatin or cholestyramine, the synergistic further increase with lovastatin and cholestyramine, and a dramatic decrease with the high cholesterol diet. Spot number 413 is the most strongly regulated protein in the present experiment, showing a 5- to 10-fold induction after a 1 week treatment with 0.075% lovastatin and 1% cholestyramine in the diet (Figs. 9 and 10). Its expression follows precisely the expectation for an enzyme whose abundance is controlled by the cholesterol level; it is progressively increased from the control levels by cholestyramine, lovastatin and lovastatin plus cholestyramine, and it sinks below the threshold of detection in animals fed the high cholesterol diet. This spot has been tentatively identified as the cytosolic HMG-CoA synthase, based on a reaction with an antiserum to that protein provided by Dr. Michael Greenspan at Merck Sharp & Dohme Research Laboratories. This enzyme lies immediately before HMG-CoA reductase in the liver cholesterol biosynthesis pathway, and is known to be co-regulated with it. Spot 413 has an SDS molecular weight of about 54 000 and a CPK pI of -11.4, in reasonably close agreement with a molecular weight of 57 300 and a CPK pI of -15.7 computed from the known sequence of the hamster enzyme [43].

Using a classical product-moment correlation test (Kepler procedure CORREL), a series of five additional spots was found to be coregulated with 413. The level of correlation was exceedingly high (> 95%). Two of these, 1250 and 933, are at similar molecular weights and approximately one charge more acidic than 413 (Fig. 9), indicating that they may be covalently modified forms of the 413 polypeptide. This suspicion is strengthened by the observation that both spots are also stained by the antibody to cytosolic HMG-CoA synthase. The remaining three correlated spots appear

to comprise an additional related pair (1253 and 1001) of around 40 kDa and a single spot (1119) of around 28 kDa. Because these two presumed proteins are present at substantially lower abundances than 413, and because the cytosolic HMG-CoA synthase is reported to consist of only one type of polypeptide, they are likely to represent other, very tightly coregulated enzymes. A second group of six spots was selected based on a regulatory pattern close to the inverse of that for spot 413 (MSN's 34, 79, 178, 182, 204, 347; data not shown). For these proteins, the lowest level of expression occurs with exposure to lovastatin plus cholestyramine and the highest level upon exposure to the high-cholesterol diet. Spots 182 and 79 are highly correlated and lie about one charge apart at the same molecular weight; they may thus be isoforms of a single protein. The other four spots probably represent additional enzymes or subunits.

### 3.3.2 MSN 235 and coregulated spots

A third group of five spots, mainly comprised of mitochondrial proteins including putative mitochondrial HMG-CoA synthase spots, showed a modest induction by lovastatin alone, but little or no effect with any of the other treatments (including the combination of lovastatin and cholestyramine; Fig. 12). This result is intriguing because lovastatin was expected to affect only the regulation of enzymes of cholesterol synthesis, which is entirely extra-mitochondrial. Three of the spots (235, 134, 144) form a closely-packed triad at approximately 30 kDa, and are likely to represent isoforms of one protein. All three spots are stained by an antibody to the mitochondrial form of HMG-CoA synthase obtained from Dr. Greenspan. Subcellular fractionation indicates a mitochondrial location. The other two spots (633 at about 38 kDa and 724 at about 69 kDa) are each present at lower abundance than the members of the triad.

### 3.3.3 An example of an anti-synergistic effect

A sixth spot (367) shows strong induction by lovastatin (two- to threefold), and about half as much induction with lovastatin plus cholestyramine, but without sharing the animal-animal heterogeneity pattern of the 235-set (Fig. 13). This protein is also mitochondrial, and represents the clearest example of an anti-synergistic effect of lovastatin and cholestyramine. The existence of such an effect demonstrates that lovastatin and cholestyramine do not act exclusively through the same regulatory pathway.

### 3.3.4 Complexity of the cholesterol synthesis pathway

Taken together, these results suggest that treatment with lovastatin alone can affect both cytosolic and mitochondrial pathways using HMG-CoA, while cholestyramine, on the other hand, either alone or in combination with lovastatin, produces a strong effect on the putative cytosolic pathway, but little or no effect on the putative mitochondrial pathway. An explanation for this difference may lie in lovastatin's effect on levels of HMG-CoA and related precursor compounds that are exchanged between the cytosol and the mitochondrion, whereas cholestyramine should affect only the cytosolic pathways directly controlled by cholesterol and bile acid levels. It remains to be explained why some

proteins of the putative mitochondrial pathway are so much more variable in their expression in all groups. An examination of all the coregulated groups suggests that quantitative statistical techniques can extract a wealth of interesting information from large sets of reproducible gels. The abundance of spots in the 413 coregulation group, for example, shows an amazing level of concordance in their relative expression among the five individuals of the lovastatin and cholestyramine treatment group. This effect is not due to differences in total protein loading, since they have already been removed by scaling, and since proteins with quite different regulation patterns can be demonstrated (e.g., Fig. 13). Such effects raise the possibility that many gene coregulation sets may be revealed through the study of a sufficiently large population of control animals (i.e., without any experimental manipulation). This approach, exploiting natural biological variation in protein expression instead of drug effects, offers an important incentive for the construction of a large library of control animal patterns.

## 4 Conclusions

Because of the widespread use of rat liver in both basic biochemistry and in toxicology, there is a long-term need for a comprehensive database of liver proteins. The rat liver master pattern presented here has proven to be an accurate representation of this system, having been matched to more than 700 gels to date. As the number of proteins identified and the number of compounds tested for gene expression effects grows, we expect this database to contribute valuable insights into gene regulation. Its practical utility in several areas of mechanistic toxicology is already being demonstrated.

Received September 11, 1991

## 5 References

- [1] O'Farrell, P. J. *Biol. Chem.* 1975, 250, 4007-4021.
- [2] Klose, J. *Humangenetik* 1975, 26, 231-243.
- [3] Scheele, G. A. *J. Biol. Chem.* 1975, 250, 5375-5385.
- [4] Iborra, G. and Buhler, J. M. *Anal. Biochem.* 1976, 74, 503-511.
- [5] Anderson, N. G. and Anderson, N. L. *Behring Inst. Mitt.* 1979, 63, 169-210.
- [6] Anderson, N. G. and Anderson, N. L. *Clin. Chem.* 1982, 28, 739-748.
- [7] Heydorn, W. E., Creed, G. J. and Jacobowitz, D. M. *J. Pharmacol. Exp. Ther.* 1984, 229, 622-628.
- [8] Anderson, N. L., Nance, S. L., Tollaksen, S. L., Gierke, F. A. and Anderson, N. G. *Electrophoresis* 1985, 6, 592-599.
- [9] Ratner, R. R. and Langley, C. H. *Biochem. Genet.* 1980, 18, 185-197.
- [10] Klose, J. *Mol. Evol.* 1982, 18, 315-318.
- [11] Neel, J. V., Baier, L., Hanash, S. J. and Erickson, R. P. *J. Hered.* 1985, 76, 314-320.
- [12] Marshall, R. R., Raj, A. S., Grant, F. J. and Heddlie, J. A. *Can. J. Genet. Cytol.* 1983, 27, 457-466.
- [13] Taylor, J., Anderson, N. G., Anderson, N. G., Gemmell, A., Giometti, C. S., Nance, S. L. and Tollaksen, S. L. in: Dunn, M. J. (Ed.), *Electrophoresis '86*. Verlag Chemie, Weinheim 1986, pp. 583-587.
- [14] Giometti, C. S., Gemmell, M. A., Nance, S. L., Tollaksen, S. L. and Taylor, J. *J. Biol. Chem.* 1987, 262, 12764-12767.
- [15] Anderson, N. L., Gierke, F. A., Nance, S. L., Gemmell, M. A., Tollaksen, S. L. and Anderson, N. G. in: Galteau, M.-M. and Sies, G. (Eds.), *Progress Recent in Electrophoresis, Bioimaging and Protein*. Universitaires de Nancy, Nancy 1986, pp. 257-260.
- [16] Anderson, N. L., Swanson, M., Gierke, F. A., Tollaksen, S. L., Gemmell, A., Nance, S. L. and Anderson, N. G. *Electrophoresis* 1986, 7, 44-48.

- Anderson, N. L., Gier, F. A., Nanc, S. L., Gemmell, M. A., Tollakson, S. L. and Anderson, N. G., *Fundam. Appl. Toxicol.* 1987, 8, 39-50.
- Anderson, N. L., in: *New Horizons in Toxicology*, Eli Lilly Symposium, 1991, in press.
- Antoine, B., Rahimi-Pour, A., Sier, G., Magdalou, J. and Galteau, M. M., *Cell. Biochem. Funct.* 1987, 3, 217-231.
- Elliott, B. M., Ramasamy, R., Stonard, M. D. and Spragg, S. P., *Biochim. Biophys. Acta* 1986, 876, 135-140.
- Huber, B. E., Heilman, C. A., Wirtz, P. J., Miller, M. J. and Thorgeirsson, S. S., *Hepatology* 1984, 4, 206-219.
- Wirth, P. J. and Vesterberg, O., *Electrophoresis* 1988, 9, 47-53.
- Witzmann, F. A. and Parker, D. N., *Toxicol. Lett.* 1991, 57, 29-36.
- Rampersaud, A., Waxman, D. J., Ryan, D. E., Levin, W. and Walt, F. G., *Jr. Arch. Biochem. Biophys.* 1985, 242, 174-183.
- Vlasuk, G. P. and Walt, F. G., *Jr. Anal. Biochem.* 1980, 105, 112-120.
- Anderson, N. G. and Anderson, N. L., *Anal. Biochem.* 1978, 85, 331-340.
- Anderson, N. L. and Anderson, N. G., *Anal. Biochem.* 1978, 85, 341-354.
- Anderson, L., Hofmann, J.-P., Anderson, E., Walker, B. and Anderson, N. G., in: Endler, A. T. and Hanzath, S. (Eds.), *Two-Dimensional Electrophoresis*, VCH Verlagsgesellschaft, Weinheim 1989, pp. 268-297.
- Anderson, L., *Two-Dimensional Electrophoresis: Operation of the ISO-DAL<sup>®</sup> System*, Large Scale Biology Press, Washington, DC 1988, ISBN 0-945532-00-8, 170pp.
- Neuhoff, V., Stamm, R. and Eibl, H., *Electrophoresis* 1985, 6, 427-448.
- [31] Neuhoff, V., Arold, N., Taube, D. and Ehrhardt, W., *Electrophoresis* 1988, 9, 255-262.
- [32] Anderson, N. L. and Hickman, B. J., *Anal. Biochem.* 1979, 92, 312-320.
- [33] Sidman, K. E., George, D. E., Barker, W. C. and Hunt, L. T., *Nucl. Acids Res.* 1988, 16, 1869-1871.
- [34] Taylor, J., Anderson, N. L., Coulter, B. P., Scandora, A. E. and Anderson, N. G., in: Radola, B. J. (Ed.), *Electrophoresis '79*, de Gruyter, Berlin 1980, pp. 329-339.
- [35] Taylor, J., Anderson, N. L. and Anderson, N. G., in: Allen, R. C. and Arnaud, P. (Eds.), *Electrophoresis '81*, de Gruyter, Berlin 1981, pp. 383-400.
- [36] Anderson, N. L., Taylor, J., Scandora, A. E., Coulter, B. P. and Anderson, N. G., *Clin. Chem.* 1981, 27, 1807-1820.
- [37] Taylor, J., Anderson, N. L., Scandora, A. E., Jr., Willard, K. E. and Anderson, N. G., *Clin. Chem.* 1982, 28, 861-866.
- [38] Taylor, J., Anderson, N. L. and Anderson, N. G., *Electrophoresis* 1983, 4, 338-345.
- [39] Anderson, N. L. and Taylor, J., in: *Proceedings of the Fourth Annual Conference and Exposition of the National Computer Graphics Association*, Chicago, June 26-30, 1983, pp. 69-76.
- [40] Anderson, N. L., Hofmann, J.-P., Gemmell, A. and Taylor, J., *Clin. Chem.* 1984, 30, 2031-2036.
- [41] Anderson, L., in: Schafer-Nielsen, C. (Ed.), *Electrophoresis '88*, VCH Verlagsgesellschaft, Weinheim 1988, pp. 313-321.
- [42] Neidhardt, F. C., Appleby, D. A., Sankar, P., Hutton, M. E. and Phillips, T. A., *Electrophoresis* 1989, 10, 116-121.
- [43] Gil, G., Goldstein, J. L., Slaughter, C. A. and Brown, M. S., *J. Biol. Chem.* 1986, 261, 3710-3716.

## 6 Addendum 1: Figures 1-13

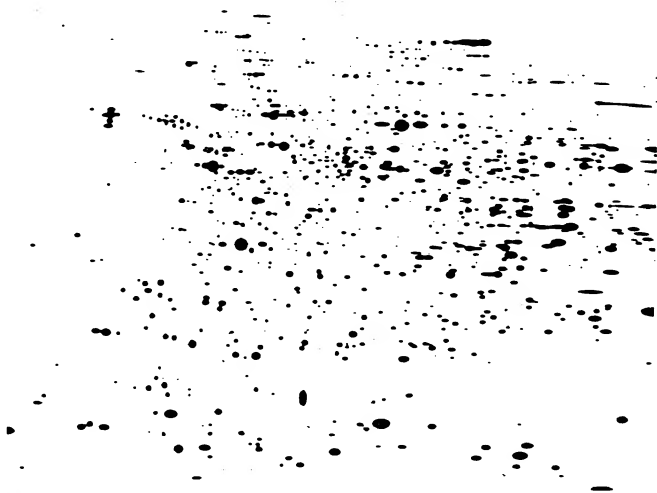


Figure 1. Synthetic representation of the standard rat liver 2-D master pattern, rendered as a greyscale image using a videoprinter.



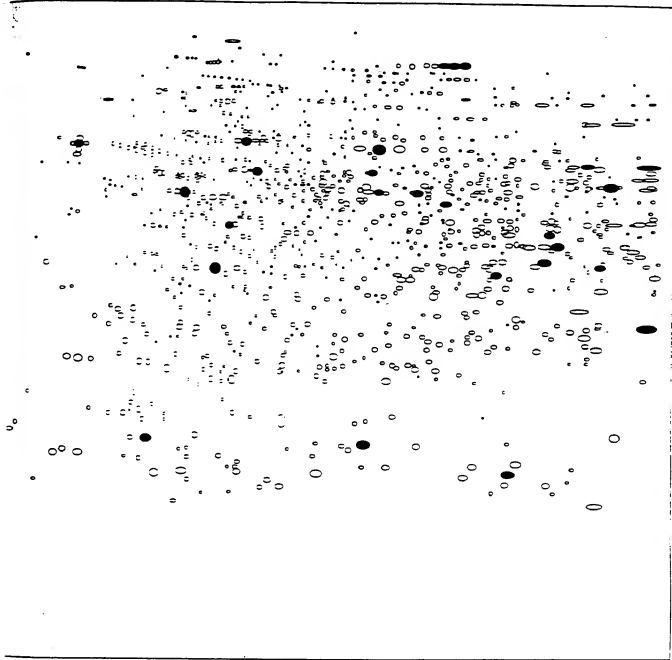


Fig. 2. Schematic representation of the master pattern (the same as Fig. 1), useful as an aid in relating specific areas of Fig. 1 and the following detailed prints.

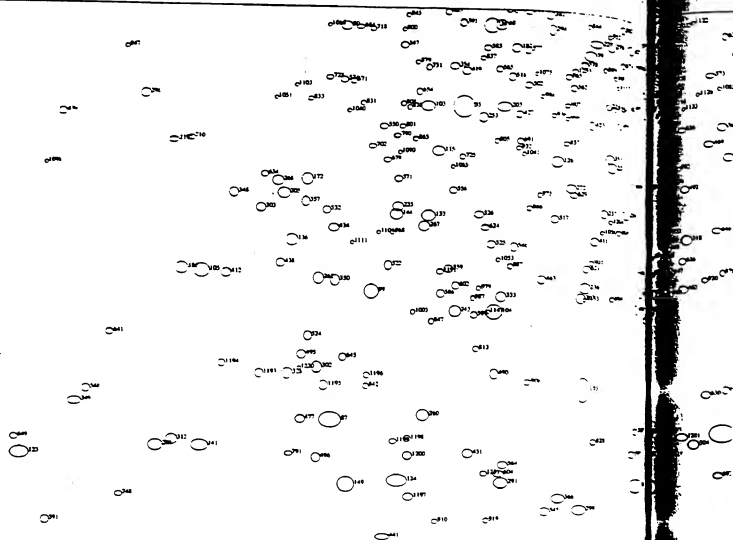
1



Figure 3. Upper left (high molecular weight, acidic) quadrant (Q1) of the rat liver map, showing spot numbers.



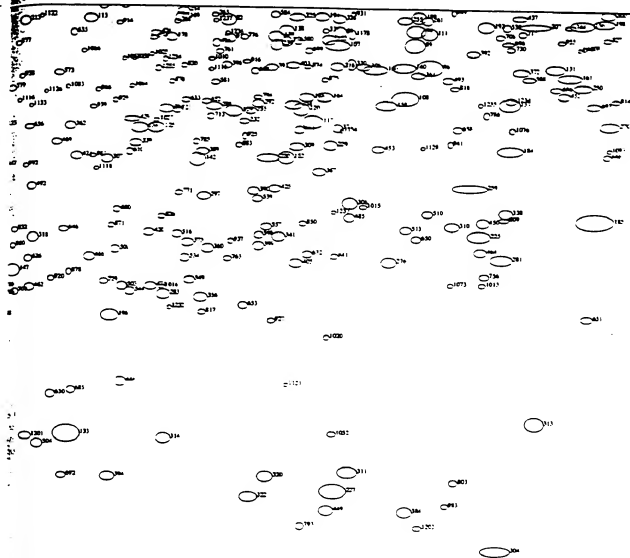
Figure 4. Upper right (high molecular weight, basic) quadrant (#2) of the rat liver map, showing spot numbers.



3

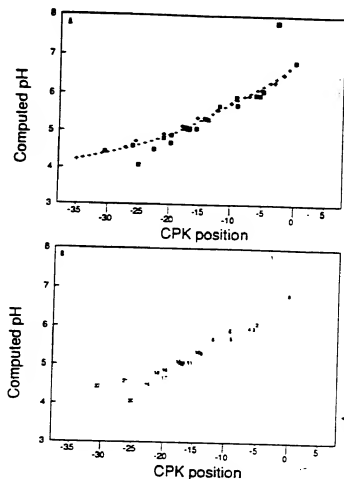
Figure 5. Lower left (low molecular weight, acidic) quadrant (#3) of the rat liver map, showing spot numbers.

4. Lower

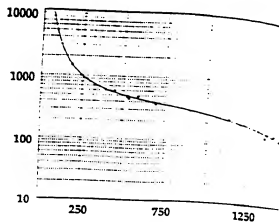


4

Fig. 6. Lower right (low molecular weight, basic) quadrant (#4) of the rat liver map, showing spot numbers.



Number of Residues



Gel Y Coordinate

Figure 8. Plot of number of amino acids versus gel Y-position, with fitted curve used to predict molecular mass of unidentified proteins.

Figure 7. (a) Plot of computed isoelectric point versus gel X-position for two sets of carbamylated standard proteins (rabbit muscle CPK [•] and human hemoglobin  $\beta$  chain, filled diamonds) and several other proteins (shaded squares). (b) The identities of the various proteins represented by the squares are indicated by the numbers in corresponding positions on (a); these refer to Table 4.

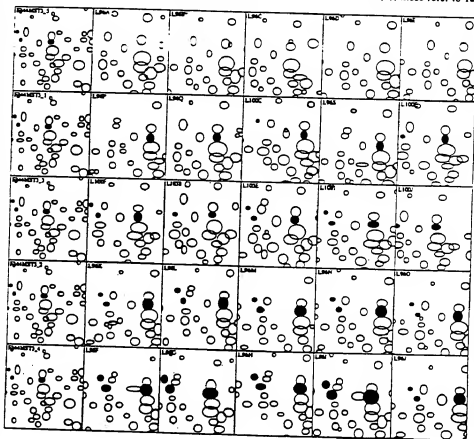


Figure 9. Montage showing effects in the region of MSN-413. The montage shows a small window into one portion of the 2-D pattern, one row of windows for each experimental group, and one panel for each gel in the experiment. The left-most pattern in each row is a group-specific copy of the master pattern followed by the patterns for the five individual rats in the group. The highlighted protein spots (filled circles) are spot 413 (on the right of each panel, identified as cytosolic HMG-CoA synthase) and two modified forms of it (1250 and 933). From the top, the rows (experimental groups) are: high cholesteryl controls, cholestyramine, lovastatin, and lovastatin plus cholestyramine.

# Regulation of Rat Liver 413

(Putative Cytosolic HMG-CoA Synthase, 53kd)  
Test Compounds in Diet

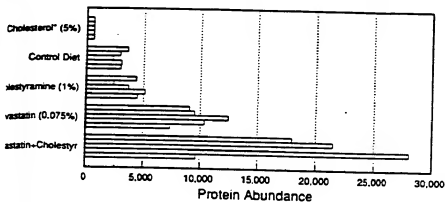


Figure 10. Bargraph showing the quantitative effects of various treatments on the abundance of MSN:413 (cytosolic HMG-CoA synthase) in the gels of Fig. 9.

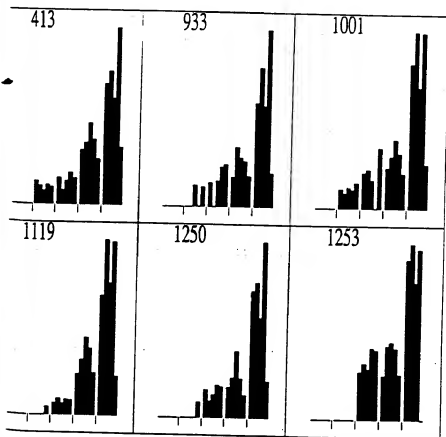


Figure 11. Bargraphs of a series of six correlated spots including MSN:413. In the bargraphs, the abundances of the appropriate spot (master spot number shown at the top of the panel) in each animal are shown. The five five-animal groups are in the order (left to right): high cholesterol, controls, cholestyramine, lovastatin, and lovastatin plus cholestyramine. Each bar within a group represents one experimental animal liver (one 2-D gel). Note the correlated expression of the 6 spots, especially in the two far right (most strongly induced) groups.

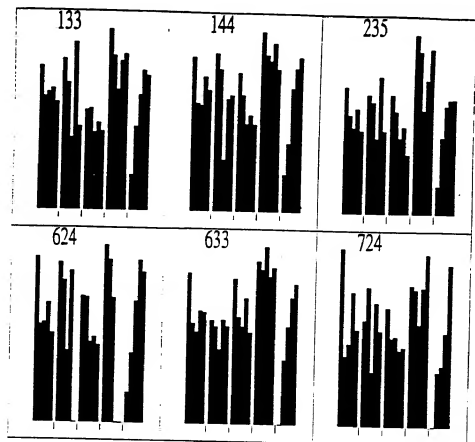


Figure 12. Data on a second coregulated group of spots, presented as in Fig. 11. The fourth experimental group (lovastatin) shows a modest induction, while the fifth group (lovastatin plus cholestyramine) does not.

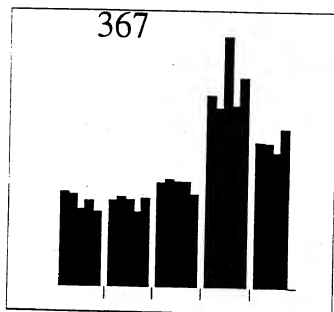


Figure 13. Data on spot MSN-367, presented as in Fig. 11. This protein shows unambiguously the anti-synergistic effect of lovastatin and cholestyramine (fifth group) as compared to lovastatin (fourth group). This response contrasts strongly with the regulation pattern seen in Fig. 11.

Appendix

Mass

X

311

548

812

845

825

835

725

646

1204

332

767

313

807

1184

1263

743

766

1215

1145

1037

863

712

763

304

1165

694

1318

1824

1252

1387

306

605

621

1113

1820

725

2001

722

678

1682

1097

1171

1400

1603

1888

735

1283

1252

779

1084

856

638

1502

1570

1284

1338

1033

1767

625

634

1811

1412

1471

1682

1817

1596

1589

516

1705

651

1415

1772

1338

1708

Further table of:  
regulated spots



Table 1. Master table of proteins in the rat liver database<sup>11</sup>

MSN	X	Y	CPKai	SOSMW	MSN	X	Y	CPKai	SOSMW	MSN	X	Y	CPKai	SOSMW
3	311	434	<-35.0	63,800	95	1119	536	-9.9	53,800	174	1364	183	-4.7	162,900
5	568	283	<-24.3	102,800	96	1731	756	-2.0	40,700	175	825	383	-15.7	69,300
8	812	426	-16.0	84,800	97	1033	566	-11.4	51,600	177	1582	553	-3.6	52,600
11	549	268	-26.2	101,000	98	1406	565	-6.1	51,700	178	1321	710	-7.2	43,000
15	845	520	-15.3	56,200	99	578	1149	-23.8	25,000	179	1088	515	-10.4	48,300
17	629	509	-21.6	50,000	100	2004	538	>0.0	53,700	180	1866	567	-0.5	51,600
18	906	414	-14.0	66,300	101	1106	623	-10.1	47,900	181	411	295	-32.1	91,200
19	755	298	-17.5	90,200	102	482	455	-28.5	61,300	182	804	730	-16.2	42,000
20	640	423	-20.9	67,900	103	665	830	-20.2	37,300	184	1860	886	-0.6	34,500
21	1204	448	-6.7	62,100	104	773	1182	-17.0	23,800	185	1997	1017	-10.0	29,800
22	332	434	<-35.0	63,800	105	312	1117	<-35.0	26,100	186	279	1113	<-35.0	50,100
23	787	424	-16.6	65,000	106	1769	509	-1.5	56,800	187	773	296	-17.0	80,800
24	513	417	<-35.0	66,000	107	1585	720	-3.6	42,500	188	1538	807	-4.2	38,400
25	807	516	-16.1	55,600	108	1692	807	-2.4	36,300	191	1650	674	-3.9	44,800
27	1184	524	-4.0	54,900	109	1482	563	-4.8	48,700	192	1818	687	-0.9	44,200
28	1263	446	-8.0	62,400	110	778	516	-16.9	55,500	193	1489	555	-5.0	52,400
29	743	605	-17.8	49,000	111	1728	700	-2.0	43,500	194	380	266	-4.4	101,600
30	768	112	-17.2	348,600	113	1191	680	-8.9	44,500	195	784	632	-16.7	47,300
32	1216	417	-4.6	66,000	114	1298	185	-7.5	160,800	196	1227	1185	-8.4	23,700
33	1145	445	-5.5	62,500	115	682	907	-19.6	34,100	197	667	553	-20.1	45,200
34	1037	555	-11.3	52,400	116	1146	610	-9.5	48,700	198	2006	681	-0.0	44,500
35	863	412	-14.9	66,600	117	1548	849	-4.1	36,500	199	1711	674	-2.2	44,800
36	712	606	-18.7	48,800	118	1050	577	-11.1	50,800	200	872	424	-14.7	65,000
38	763	694	-17.3	43,800	120	1530	828	-4.3	37,400	201	262	435	<-35.0	63,700
39	304	470	<-35.0	58,800	121	838	423	-11.4	65,200	202	736	553	-18.0	107,800
41	1165	569	-9.2	51,400	122	1572	712	-3.8	42,900	203	786	829	-16.7	37,400
42	684	607	-19.6	48,800	123	23	1433	<-35.0	15,300	204	1224	589	-5.5	50,000
43	1318	589	-7.3	50,000	124	621	1474	-21.9	13,900	205	439	983	-30.9	31,100
44	1924	365	-0.1	74,600	125	1298	862	-7.1	36,000	206	1994	1017	-10.0	51,300
46	1203	586	-8.7	50,200	126	872	921	-14.7	33,500	207	1895	687	-0.3	36,000
47	1391	447	-6.3	62,300	127	1000	717	-12.0	42,600	208	240	1418	<-35.0	15,800
48	309	454	<-35.0	61,500	128	1229	311	-8.4	86,100	210	1700	499	-2.3	57,000
49	605	587	-22.5	50,100	129	1422	832	-5.8	37,300	211	902	517	-14.1	55,400
50	621	535	-21.8	53,900	130	1776	499	-1.4	57,000	212	884	1014	-10.4	44,400
51	1113	522	-10.0	55,000	131	1930	757	-0.1	40,700	214	1340	668	-7.0	45,200
52	1820	499	-0.9	57,000	132	660	537	-20.4	53,800	215	1591	495	-3.5	57,300
53	725	177	-18.3	170,800	133	666	1019	-20.2	29,700	216	1585	755	-3.6	40,700
54	2001	500	>0.0	56,900	134	1271	862	-7.9	36,000	217	1159	393	-9.3	69,300
55	722	830	-18.4	37,300	135	1161	1369	-9.3	16,800	218	831	572	-13.5	51,200
56	778	533	-19.8	54,100	136	453	1063	-29.7	28,100	219	713	177	-18.7	170,500
57	1682	302	-2.5	89,000	137	1858	823	-0.6	37,700	220	1479	911	-4.9	33,900
58	1091	580	-10.3	50,500	138	1504	697	-4.6	43,700	221	965	927	-12.8	33,300
59	1171	585	-9.2	50,300	139	1488	707	-4.8	43,200	223	834	716	-13.5	42,700
60	1400	624	-6.2	47,800	140	1699	756	-2.4	40,700	225	1812	1045	-1.0	28,800
61	1853	508	-0.6	56,200	141	311	1417	<-35.0	15,800	226	821	411	-15.8	66,800
62	1888	567	-0.4	51,500	142	1366	915	-5.7	33,800	227	1586	1483	-3.6	13,600
65	735	297	-18.1	90,500	143	1429	346	-5.7	77,900	228	1065	567	-10.8	51,600
66	1263	312	-6.0	85,900	144	615	1017	-22.1	29,800	229	1577	890	-3.7	34,800
67	1252	407	-8.1	67,300	145	2006	566	>0.0	51,600	230	1458	496	-5.2	57,300
68	779	692	-16.8	43,900	146	2006	518	>0.0	55,300	232	1440	849	-5.5	36,500
69	1064	286	-10.8	90,800	147	1070	1108	-10.7	26,500	234	1692	489	-2.4	57,900
71	656	589	-20.6	50,000	148	1347	578	-6.9	50,800	235	618	1004	-22.0	30,300
72	638	545	-21.1	49,100	149	541	1481	-25.7	13,700	236	920	1138	-13.7	25,400
73	1582	583	-3.6	50,400	150	1645	760	-2.8	40,500	237	952	1008	-13.2	38,200
74	1570	556	-3.8	52,300	151	1269	236	-7.9	117,000	238	1611	541	-3.2	53,500
75	1264	621	-8.0	48,000	152	1507	911	-4.5	33,900	239	1489	720	-4.8	42,500
76	1338	564	-7.0	51,800	153	1722	448	-2.1	62,100	240	501	448	-27.7	62,100
77	1833	363	-0.8	54,400	154	932	523	-13.5	56,500	241	1820	569	-0.9	51,400
78	1767	565	-1.5	51,700	155	1031	294	-11.4	91,400	242	1357	658	-6.8	45,800
79	925	738	-13.6	41,600	156	1970	684	>0.0	44,400	243	711	1182	-18.7	23,800
80	534	698	-26.1	43,600	157	1258	183	-8.1	162,400	244	1855	621	-0.6	48,000
81	1811	363	-1.0	74,500	158	1275	417	-7.8	65,900	245	1189	474	-8.9	59,300
82	1412	681	-6.0	44,500	159	1863	820	-2.5	37,800	246	551	456	-25.1	61,000
83	1471	347	-5.0	77,500	160	1034	527	-11.4	54,600	247	1348	604	-4.9	49,100
84	1662	563	-2.7	51,800	161	1953	771	>0.0	40,000	248	460	448	-29.3	62,100
85	1596	479	-3.4	36,900	162	1020	1482	-11.6	13,700	249	1733	451	-1.9	61,800
86	1817	301	-8.9	89,100	163	1566	806	-3.8	38,400	250	1874	786	-0.0	39,200
87	516	1371	-27.0	17,400	164	1905	565	-11.2	51,700	251	808	392	-16.1	61,400
88	1598	698	-3.5	43,600	167	1340	181	-7.0	164,900	252	874	553	-14.6	52,500
89	1706	719	-2.2	42,500	168	1506	583	-4.6	50,400	253	753	848	-17.6	36,500
90	651	329	-20.8	81,700	169	1338	678	-7.0	44,700	254	995	450	-12.1	61,900
91	1415	710	-6.0	43,000	170	1869	541	>0.0	53,500	255	1690	679	-2.4	44,600
92	1773	545	-1.4	53,200	171	800	378	-16.3	71,800	256	984	1006	-12.1	30,400
93	1338	448	-7.0	62,300	172	476	958	-28.7	32,100	257	508	464	-27.4	60,400
94	1708	696	-2.2	43,700	173	919	1314	-13.7	19,300	258	1517	820	-4.4	37,800

Master table of proteins in the rat liver database, showing spot master number, gel position (x and y), isoelectric point relative to CPK standards, and predicted molecular mass (from the standard curve of Fig. 8).

MSN	X	Y	CPKd	SDSMW
250	1796	961	-1.1	31,800
260	661	1361	-20.4	17,700
261	1725	679	-2.0	44,600
262	486	1127	-28.0	25,800
263	1063	172	-10.9	177,400
265	1390	673	-6.3	45,000
266	510	437	-27.3	26,000
267	660	1038	-20.4	63,400
268	430	961	-31.0	31,800
269	1044	806	-11.2	48,900
270	2019	652	-15.0	36,300
271	857	422	-6.5	65,200
272	895	968	-14.2	31,700
274	1292	712	-7.6	42,900
275	1350	560	-6.9	49,800
276	1670	1089	-2.6	27,100
277	698	538	-18.4	53,700
278	961	718	-13.0	42,600
279	879	570	-14.5	51,300
281	1848	1064	-0.7	27,300
282	1505	525	-4.8	54,600
283	1313	1147	-7.3	25,100
284	1314	829	-7.3	37,400
285	1332	408	-7.1	67,200
286	1277	852	-7.8	46,100
288	138	89	-8.5	37,800
289	1147	579	-9.5	50,700
290	925	511	-13.6	55,800
291	787	1476	-16.6	13,800
292	1462	818	-5.1	37,800
293	423	449	-25.3	62,000
294	600	698	-14.9	43,600
295	1162	609	-3.3	48,700
296	218	814	-35.0	38,000
297	1377	979	-6.5	31,300
299	913	1523	-13.9	12,400
300	2012	667	-30.0	45,300
301	702	176	-19.0	188,200
302	494	1280	-28.1	20,400
303	403	1008	-32.6	30,100
304	1843	1585	-10.7	49,800
305	1040	563	-11.1	10,300
306	1608	989	-3.3	30,900
307	1219	916	-8.5	33,700
308	1527	795	-3.0	40,700
309	1524	802	-4.4	34,700
310	1789	1028	-1.5	29,400
311	1609	1451	-3.3	14,700
312	266	1408	-35.0	20,100
313	1902	1365	-0.3	17,600
314	1316	1395	-7.3	16,600
315	1341	529	-7.0	54,900
318	1140	1053	-10.1	28,500
320	1480	1439	-4.8	14,400
321	850	620	-15.1	49,100
322	1454	1464	-5.3	13,300
323	707	626	-20.0	47,700
324	855	101	-20.6	420,500
325	1321	675	-4.4	44,600
326	1587	677	-3.8	44,700
327	1388	620	-6.3	67,000
328	448	1291	-30.0	67,000
330	1608	751	-3.3	40,800
331	1566	697	-3.8	43,700
332	521	471	-26.3	59,600
333	784	1156	-16.7	59,600
334	1059	407	-10.9	67,300
335	1593	303	-3.5	88,500
336	1618	586	-3.2	49,400
338	1854	1004	-0.6	34,800
339	1265	888	-8.0	30,300
340	581	565	-23.6	50,300
341	1497	1047	-4.7	26,700
343	1351	295	-6.8	102,200
344	1813	549	-0.9	52,800

MSN	X	Y	CPKd	SDSMW
345	1006	578	-11.9	50,800
346	1095	640	-10.3	46,800
347	625	728	-21.7	42,000
348	961	983	-35.3	31,100
349	110	1343	-35.0	18,300
350	521	1130	-26.7	25,700
351	912	619	-13.9	48,100
352	1574	530	-3.7	54,300
353	961	762	-18.9	40,400
354	706	912	-12.9	33,800
355	1450	830	-5.3	37,300
356	1374	1152	-6.5	24,900
357	474	997	-28.7	30,600
358	784	346	-16.3	77,800
359	784	338	-6.4	27,900
360	1364	1068	-6.4	27,900
361	1713	789	-2.1	40,100
362	1161	859	-9.3	36,100
363	614	1156	-13.8	24,800
364	412	435	-32.0	79,400
365	741	486	-17.9	56,700
366	78	1503	-14.6	13,000
367	1560	635	-3.9	33,000
368	963	620	-12.4	55,200
369	434	611	-31.0	63,000
370	639	610	-21.2	48,700
371	1587	860	-3.6	36,100
372	1875	762	-0.5	40,400
373	1313	1058	-4.6	28,300
374	1506	715	-8.8	42,700
375	1823	532	-0.9	54,200
376	254	417	-35.0	65,900
377	1409	583	-6.1	50,400
378	521	594	-21.8	57,500
379	1017	595	-11.7	49,600
381	953	586	-13.1	49,400
382	856	674	-15.0	44,800
383	1252	258	-8.1	105,300
384	1699	1518	-2.3	12,500
385	1042	453	-11.2	57,500
386	1490	583	-4.7	50,400
387	1554	603	-4.0	49,100
388	1193	604	-8.9	67,700
389	1374	902	-6.5	34,300
390	1456	969	-5.2	31,700
391	718	690	-18.5	44,000
392	1799	732	-1.1	41,900
393	1482	758	-4.8	40,600
394	1227	1461	-8.4	14,400
395	1530	577	-4.3	50,800
396	1410	755	-6.0	40,800
397	912	256	-13.9	106,400
398	1465	1063	-4.9	61,900
400	1473	450	-5.0	29,100
401	1029	1140	-11.5	25,300
403	1516	754	-4.4	40,800
404	1495	554	-4.7	52,500
405	1525	1002	-4.9	27,100
406	723	252	-18.8	106,000
408	650	663	-20.8	45,500
410	1501	478	-6.6	59,000
411	936	1057	-13.4	28,300
412	350	1120	-35.9	26,000
413	1033	538	-11.4	53,700
415	737	425	-18.0	64,900
416	1578	606	-3.7	48,900
417	646	496	-21.0	57,300
418	1695	862	-2.3	58,600
419	725	770	-18.3	50,400
420	1289	1041	-7.7	28,900
421	1171	912	-9.1	33,900
422	582	162	-22.8	163,700
423	929	856	-13.6	36,200
424	739	625	-17.9	47,700
425	1490	965	-4.7	31,800

MSN	X	Y	CPKd	SDSMW
426	1296	704	-7.8	43,300
427	810	843	-16.0	36,800
428	1565	303	-3.9	36,800
429	1259	847	-8.0	88,700
430	1253	562	-8.1	36,800
431	734	1426	-18.1	51,800
432	483	1041	-11.6	51,800
434	518	1041	-33.5	63,800
435	1020	1170	-28.9	28,800
436	1122	196	-9.8	24,300
437	1670	673	-0.5	147,800
438	435	1102	-31.0	45,000
439	86	847	-35.0	26,700
440	1740	544	-1.8	36,800
441	599	1571	-22.8	53,200
443	743	335	-17.8	10,800
446	801	668	-16.2	80,100
447	1050	926	-11.1	45,200
448	1245	1290	-8.2	33,300
449	1576	1516	-3.7	19,800
450	1094	440	-0.0	34,800
451	1818	1021	-0.9	26,800
452	1045	802	-10.3	63,100
453	1652	864	-2.8	38,800
454	1403	500	-6.1	56,900
456	1384	718	-6.3	42,600
457	905	905	-11.1	63,500
458	1038	581	-11.1	50,600
460	1598	294	-3.4	91,400
461	1528	963	-4.3	35,900
462	1098	1117	-10.2	25,400
463	849	1072	-15.2	25,800
464	1814	1072	-0.8	80,300
465	1388	481	-6.3	56,700
466	1184	1064	-8.9	27,300
468	577	467	-23.9	60,100
469	1140	88	-9.6	34,900
470	797	524	-1.1	54,800
471	1293	1133	-7.6	25,900
472	618	655	-21.9	46,000
473	2009	299	-0.0	89,900
474	1205	215	-8.7	121,300
475	1035	788	-11.4	39,300
476	160	155	-35.0	207,800
477	469	1370	-28.9	17,400
478	599	662	-22.8	45,800
479	1009	540	-11.8	53,900
480	1216	235	-8.6	117,400
482	816	346	-15.9	77,800
483	683	673	-19.3	44,800
485	1608	4313	-3.3	30,000
486	478	599	-26.6	49,300
487	1025	607	-11.5	48,800
488	1045	1186	-11.2	23,700
489	1609	301	-3.3	89,200
490	775	1299	-17.0	20,100
491	692	178	-19.3	169,300
492	1100	964	-10.2	31,800
493	1760	776	-1.6	39,700
494	882	247	-14.5	50,900
495	407	1258	-28.9	21,200
496	484	1436	-28.1	15,200
497	980	852	-12.5	36,400
499	1414	546	-6.0	53,900
500	1234	1072	-6.3	27,800
501	1246	659	-8.2	45,700
502	824	732	-15.7	39,000
503	1246	1134	-8.2	25,300
504	1115	1407	-9.9	15,200
505	1189	391	-8.9	68,700
506	1578	402	-3.7	68,000
507	787	250	-16.6	108,000
508	979	552	-12.5	52,600
509	1153	619	-9.4	48,100
510	1730	1006	-2.0	30,700

MSN	X
800	
1009	
1606	
948	
461	
1334	
868	
798	
622	
147	
1332	
503	
1190	
479	
768	
747	
1170	
1502	
1728	
306	
507	
870	
1347	
1513	
803	
11851	
1463	
809	
1164	
803	
1259	
803	
1182	
1355	
1300	
1126	
1365	
996	
1300	
902	
700	
1028	
788	
880	
1212	
780	
1142	
532	
771	
1048	

SEN	X	Y	CPKpl	SDSMW
511	800	484	-16.0	56,400
512	1009	533	-10.2	54,100
513	1696	1034	-2.3	29,200
514	948	636	-13.2	47,100
515	481	543	-28.5	53,400
516	1334	1044	-7.1	28,800
517	856	1021	-14.8	29,700
518	798	779	-11.3	39,600
519	822	670	-17.7	47,100
520	632	185	-21.5	186,000
521	1332	830	-7.1	37,300
522	603	1104	-22.6	26,600
523	1190	309	-8.9	86,800
524	479	1226	-26.6	22,300
525	768	1066	-17.2	28,000
526	747	1016	-17.7	29,800
527	1170	1231	-4.6	53,400
528	1502	542	-9.2	119,600
530	1726	620	-2.0	48,000
532	507	1011	-27.4	30,000
533	870	489	-17.7	57,900
534	1347	1085	-6.9	27,300
535	1513	346	-8.5	77,800
536	306	654	-35.0	46,000
538	1651	689	-0.7	44,100
539	1463	982	-5.1	31,100
540	909	561	-13.8	52,000
541	625	289	-21.7	83,100
542	1164	198	-9.2	146,200
543	803	655	-16.2	45,900
544	1259	1143	-8.0	25,200
545	856	1526	-15.0	57,800
546	803	1071	-16.2	12,200
547	1162	274	-8.3	96,400
548	1226	1321	-35.0	19,000
549	1355	1122	-4.6	25,900
550	595	866	-22.0	35,800
552	1369	494	-6.6	57,500
553	992	405	-12.2	67,600
555	1125	410	-9.6	66,900
556	705	875	-18.9	31,400
557	1477	1030	-4.9	29,300
558	980	583	-12.5	50,400
559	700	1109	-19.1	26,400
560	1028	621	-11.5	48,000
562	898	794	-14.1	38,900
563	798	1446	-16.6	14,900
565	777	766	-16.9	40,200
566	980	328	-12.5	61,900
567	1518	611	-4.4	46,600
568	1212	601	-8.6	45,600
569	760	584	-17.4	48,700
571	618	956	-21.9	32,100
573	1142	771	-9.6	40,000
574	532	787	-26.2	39,300
575	771	250	-17.1	109,200
576	1068	534	-10.8	54,100
577	822	734	-15.7	41,800
578	914	754	-13.8	40,800
579	1064	794	-10.8	38,900
580	1524	714	-4.4	42,800
581	1392	783	-6.3	39,400
582	982	686	-12.4	44,200
584	1487	672	-4.8	45,000
586	758	731	-17.4	41,900
588	697	1152	-19.5	24,900
589	930	523	-13.5	56,000
588	1888	774	-0.4	39,300
590	642	485	-21.1	56,300
591	1317	519	-7.2	55,300
592	1014	814	-11.7	51,500
593	732	176	-18.1	172,300
594	1627	478	-3.0	59,000
595	1009	1426	-11.6	15,500

MSN	X	Y	CPKpl	SDSMW
596	619	269	-21.9	100,500
597	1176	461	-9.1	60,700
598	1465	1044	-5.0	28,800
599	741	1188	-17.9	23,600
600	807	403	-14.0	66,000
601	687	658	-19.5	45,800
602	712	1138	-18.7	25,400
603	898	181	-14.1	185,200
604	783	1461	-16.7	14,400
605	736	223	-16.0	125,300
606	629	273	-21.6	98,700
607	1064	286	-10.6	94,000
608	883	503	-14.5	56,700
609	1212	610	-30.0	48,700
610	1255	903	-8.1	34,200
612	1103	391	-10.1	69,600
613	778	265	-16.9	102,000
614	824	518	-15.7	55,400
615	1005	195	-10.3	149,100
616	1759	479	-1.6	59,000
617	994	372	-12.1	72,900
618	751	374	-17.6	72,400
619	1429	516	-5.7	55,300
620	1050	520	-11.1	55,200
621	923	1105	-13.7	26,600
622	1462	622	-5.1	47,900
623	759	225	-17.4	124,000
624	758	1038	-17.4	29,000
625	1436	636	-5.5	48,900
626	1096	1089	-10.2	27,200
627	942	548	-13.3	53,000
628	899	621	-16.0	48,000
629	899	979	-14.1	31,300
630	1135	1321	-8.6	19,100
631	979	615	-12.5	71,300
632	1542	1076	-4.1	27,600
633	1345	814	-6.9	38,000
634	409	950	-32.2	32,400
635	1165	704	-8.2	43,300
636	774	604	-17.0	49,000
637	1263	524	-8.0	54,800
638	952	411	-13.1	66,700
639	1717	575	-2.1	51,000
640	994	292	-12.1	52,000
641	165	1224	-35.0	22,400
642	803	251	-16.2	108,900
643	719	296	-18.5	90,700
644	1004	294	-10.2	91,400
645	534	1263	-26.1	20,000
646	1153	1038	-9.4	29,000
648	1246	204	-8.2	140,000
649	1743	1406	-35.0	16,200
650	1048	1048	-2.1	26,600
651	1986	1183	-0.0	23,800
652	1378	616	-6.5	38,200
653	1442	1165	-5.5	24,400
654	650	806	-20.8	38,400
655	1111	551	-10.0	52,700
656	1095	861	-10.3	36,000
657	1524	540	-4.4	53,600
658	1777	860	-1.4	33,000
659	391	584	-33.4	50,400
660	977	565	-12.5	51,700
661	658	166	-20.5	67,500
662	732	312	-18.1	86,100
663	1787	567	-1.2	51,500
664	888	268	-14.4	100,900
665	899	711	-14.3	39,800
666	715	221	-18.6	26,500
667	781	227	-16.8	122,400
668	646	185	-20.1	189,100
669	1116	353	-9.9	76,300
670	1382	643	-6.4	46,800
671	547	789	-25.3	39,200
673	964	746	-12.4	41,200

MSN	X	Y	CPKpl	SDSMW
674	1661	448	-2.7	62,100
675	1523	562	-4.4	51,800
676	708	642	-18.8	46,700
677	919	615	-13.7	48,300
678	1085	551	-10.5	52,700
679	600	623	-22.7	33,400
680	1227	1004	-8.3	30,300
681	1013	283	-10.1	95,100
682	1406	477	-6.1	56,100
683	1596	249	-3.4	109,800
684	555	689	-24.8	43,500
685	1167	1313	-9.2	19,300
686	1932	790	0.0	39,100
687	1545	619	-4.1	48,100
688	1456	764	-5.2	40,300
689	1011	953	-11.8	32,300
690	1995	270	-30.0	100,200
691	812	886	-16.0	34,900
692	1154	1461	-9.4	14,400
693	1993	619	-30.0	37,800
694	228	658	-11.8	45,900
695	928	254	-7.1	107,000
696	1854	715	-0.6	42,700
697	1997	345	-30.0	78,000
698	957	563	-13.0	51,800
699	1540	730	-4.2	42,000
700	577	900	-22.8	34,400
703	1610	562	-3.2	51,900
705	1278	571	-7.8	51,200
706	1841	704	-0.7	43,300
707	1018	1386	-11.7	16,900
709	1074	1145	-10.7	25,100
710	293	889	-35.0	34,800
712	720	412	-18.5	66,600
713	1386	841	-6.4	36,800
714	1328	263	-7.1	103,100
715	698	433	-19.1	53,900
716	701	481	-19.0	58,700
717	1875	699	-0.5	43,600
718	575	702	-23.9	43,400
719	1216	204	-8.6	140,400
721	1069	464	-10.8	50,400
722	1272	506	-7.9	56,400
723	958	822	-13.0	37,700
724	763	395	-17.3	69,100
725	720	916	-18.5	33,700
726	1476	415	-4.9	66,200
727	1846	473	-0.7	59,400
728	510	783	-27.3	39,400
729	1217	1126	-8.6	25,800
730	1856	724	-0.6	42,300
731	665	765	-20.2	40,300
733	1321	312	-7.2	85,900
734	719	427	-18.5	64,600
735	1101	473	-10.2	59,500
736	1359	569	-6.7	51,000
738	696	220	-19.2	127,600
739	687	409	-19.5	67,000
740	1205	556	-8.7	106,200
741	995	563	-12.1	51,800
742	998	596	-14.1	80,500
743	881	181	-14.5	165,900
744	1951	686	-30.0	44,200
745	726	168	-16.3	183,600
746	999	643	-12.0	46,500
748	182	1503	-35.0	13,000
749	2005	649	-30.0	46,300
750	1448	575	-5.4	51,000
751	792	266	-16.5	101,900
752	465	206	-28.9	80,500
754	664	254	-20.3	107,000
755	1195	184	-8.8	161,000
756	1821	1113	-0.9	26,300
757	909	246	-13.9	111,000
760	790	133	-16.5	264,900

MSN	X	Y	CPKd	SDSMW
761	1399	733	-6.2	41,800
763	1416	1085	-5.9	27,300
764	2020	590	>0.0	51,400
765	851	475	-20.8	59,200
766	1052	1149	-11.1	25,000
767	1968	486	>0.0	59,900
768	1330	685	-7.1	44,300
769	1870	813	>0.0	48,500
770	857	617	-15.0	48,200
771	1337	974	-7.0	31,500
772	1576	502	-3.7	56,700
775	969	824	-12.8	37,600
775	1438	708	-5.5	43,100
777	1539	458	-4.2	63,400
778	850	434	-15.1	63,800
779	700	411	-19.1	66,800
780	1052	1136	-11.1	25,500
784	1413	529	-6.0	54,400
785	1364	585	-6.7	35,000
786	1822	835	-0.9	37,100
787	893	392	-14.3	69,500
790	816	882	-22.0	35,100
791	451	1429	-29.8	15,400
792	777	777	-15.9	72,000
794	1536	1543	-4.2	11,700
794	1461	807	-5.1	38,300
796	388	546	-33.6	53,100
797	1126	212	-9.8	133,700
798	933	437	-11.7	83,400
799	1420	593	-5.9	49,800
800	1759	279	-1.6	96,500
801	624	865	-21.7	35,800
802	986	522	-14.2	53,000
803	1775	1468	-1.4	44,100
804	573	196	-24.0	148,400
805	203	494	-35.0	57,400
806	980	1039	-12.5	29,000
807	902	306	-14.1	87,200
808	625	827	-21.7	37,500
809	1851	1015	-0.7	29,900
810	440	573	-30.9	51,100
811	1358	249	-6.8	109,700
812	851	383	-15.1	69,400
813	745	1246	-17.8	21,600
814	2028	810	>0.0	38,200
815	1086	645	-10.4	46,500
816	629	313	-21.6	85,700
817	1375	1177	-45.5	24,000
818	1771	790	-1.4	39,100
819	1045	263	-11.2	103,100
820	984	362	-12.4	74,600
821	1712	278	-2.2	96,700
822	1256	205	-2.1	139,200
823	1517	654	-4.4	46,000
824	1442	449	-5.5	62,000
825	1240	513	-8.3	55,800
826	1309	1014	-7.4	29,900
827	2012	700	>0.0	43,100
828	837	1405	-13.4	16,200
830	1342	756	-7.0	40,700
831	562	826	-24.5	37,800
832	1073	1039	-10.7	29,000
833	481	820	-28.5	37,800
834	501	581	-27.8	50,500
837	751	748	-17.6	41,100
838	635	833	-21.3	37,200
839	1494	459	-4.7	80,300
840	1952	301	>0.0	60,900
841	585	1080	-3.6	27,500
842	571	1312	-24.1	19,400
843	1325	649	-7.2	46,300
844	1727	301	-21.5	44,600
845	630	879	-2.0	89,200
846	2016	905	>0.0	34,200
847	673	1200	-19.9	23,200

MSN	X	Y	CPKd	SDSMW
848	1863	271	-0.5	99,500
849	1166	523	-9.2	54,900
850	1535	1024	-4.2	29,600
851	1035	826	-11.4	37,500
852	834	542	-15.5	53,400
855	499	201	-27.8	127,100
856	1053	194	-10.9	150,500
857	887	890	-14.4	34,800
858	1448	639	-5.4	45,900
859	705	311	-18.9	86,200
860	1070	1066	-10.7	75,600
861	472	447	-26.8	77,600
862	674	480	-19.9	58,800
864	1307	499	-7.4	57,000
865	645	887	-21.0	34,900
866	827	1004	-15.6	30,300
868	685	494	-19.5	57,400
869	1807	402	-1.0	68,000
870	1323	783	-7.2	39,400
871	1228	911	-8.4	29,300
872	1904	346	-0.3	77,700
873	556	647	-24.8	45,400
874	1540	756	-4.2	40,700
875	1566	777	-3.8	39,700
876	1196	351	-8.8	76,800
877	1076	720	-10.6	42,500
878	1161	1111	-9.3	26,400
879	647	757	-20.9	40,700
880	1758	594	-1.6	49,700
881	1514	278	-1.7	87,100
883	1432	890	-5.7	34,800
884	922	689	-13.7	44,100
885	1103	414	-10.1	66,400
887	1501	607	-4.6	48,900
888	798	1031	-16.3	25,800
888	636	634	-21.3	47,200
889	951	759	-13.1	40,600
890	717	548	-18.6	52,900
891	1123	229	-9.8	121,200
892	891	413	-14.3	66,400
894	1245	334	-8.2	117,800
895	1962	346	>0.0	77,700
896	1322	626	-7.2	47,700
897	420	570	-31.4	51,300
898	652	428	-20.3	64,500
899	845	243	-15.3	113,000
900	624	703	-21.7	43,400
901	931	1094	-13.5	27,000
903	799	229	-16.3	121,000
904	765	520	-17.2	55,200
905	775	889	-17.0	34,800
907	888	824	-14.4	37,600
908	828	1303	-15.6	19,700
910	581	1544	-19.7	89,100
911	1544	301	-4.1	89,100
913	1606	387	-3.3	70,400
914	1237	688	-8.3	41,100
916	1442	749	-5.5	41,100
917	1286	367	-8.0	73,700
919	764	1541	-17.3	11,700
920	1133	1123	-9.7	25,900
921	1123	380	-9.8	71,500
923	829	242	-15.6	113,200
924	1131	318	-9.7	84,300
925	1441	874	-5.5	35,400
926	679	219	-19.7	128,200
927	1487	1191	-4.8	23,500
928	1082	775	-10.5	39,800
929	1221	816	-4.4	38,000
931	1609	670	-3.3	45,100
932	810	900	-16.0	34,400
933	965	520	-12.8	55,100
934	647	462	-13.2	60,600
935	843	148	-14.8	36,800
937	1421	1056	-5.9	28,400

MSN	X	Y	CPKd	SDSMW
938	1187	827	-8.8	37,500
941	1765	885	-1.5	35,000
942	602	472	-22.7	59,600
943	312	498	-35.0	59,600
944	993	491	-12.1	57,100
945	1300	269	-7.5	100,300
946	633	429	-35.0	65,100
947	187	736	-21.6	65,100
948	1380	344	-7.6	41,800
949	1766	665	-1.5	78,300
950	1038	193	-11.3	151,000
951	860	152	-13.0	213,000
952	957	701	-14.9	213,000
954	503	547	-27.6	43,400
955	1038	712	>0.0	42,800
957	1010	816	-11.8	37,800
959	768	174	-15.3	174,900
960	596	419	-23.0	67,600
961	557	409	-24.8	65,700
962	887	320	-14.4	83,800
963	964	334	-24.5	80,500
964	969	1156	-12.8	24,800
965	671	255	-20.0	100,600
966	1204	798	-6.7	38,700
967	910	154	-13.9	210,300
968	609	1048	-22.3	29,700
969	1265	206	-7.7	71,700
970	822	232	-15.8	138,900
971	976	437	-12.6	63,400
972	403	567	-32.6	51,600
973	844	995	-35.0	57,400
974	1224	295	-9.8	91,400
977	994	664	-12.1	45,400
978	1512	642	-3.2	46,700
979	749	611	-17.9	25,300
980	1064	642	-10.8	46,700
981	1197	911	-8.8	33,900
983	1762	1508	-1.6	12,800
984	1344	317	-6.9	84,700
985	1024	1105	-11.5	26,000
987	739	1159	-17.9	24,600
988	816	555	-15.9	52,400
990	785	361	-16.7	74,900
991	1159	317	-9.3	84,500
992	1090	928	-10.4	33,300
993	1030	701	-11.5	43,400
994	847	811	-15.2	38,200
995	902	461	-14.1	60,700
996	888	847	-14.4	36,600
997	1815	579	-0.9	50,700
998	1205	504	-8.7	56,500
999	617	289	-22.0	83,100
1000	968	290	-12.8	92,700
1001	970	717	-12.4	92,700
1002	1736	479	-1.9	98,700
1003	643	1184	-21.1	23,700
1006	822	487	-15.8	58,100
1007	875	279	-14.6	66,600
1009	291	644	-35.0	46,400
1010	1386	745	-6.4	41,200
1011	459	541	-29.4	53,500
1012	679	661	-19.7	45,600
1013	1818	1128	-0.9	25,800
1014	1032	634	-11.4	47,200
1015	1629	994	3.0	30,700
1016	311	1134	-7.4	25,500
1017	1722	424	-2.0	65,000
1018	1015	743	-11.7	41,200
1020	1512	1219	-2.2	22,500
1021	781	464	-16.8	58,400
1022	1129	83	-9.7	591,000
1023	812	317	-15.9	84,800
1024	785	446	-16.7	62,400
1025	1290	739	-7.7	41,900

425  
1296  
856  
1264  
945  
1547  
1381  
1525  
1229  
1226  
1761  
541  
818  
940  
1439  
1540  
1576  
1086  
940  
426  
1563  
779  
1613  
1380  
284  
1261  
393  
1817  
1245  
1258  
705  
1181  
806  
508  
873  
1768  
806  
1963  
826  
971  
1546  
1157  
620  
1867  
2019  
1546  
61  
1954  
588  
457  
1884  
1714  
1717  
1676  
547  
1348  
1385  
1078  
975  
1202  
1022  
1905  
1211  
1114  
1464  
1048  
1122  
1022  
1006  
1830  
764  
1968

MSN	X	Y	CPKpt	SDS/MW	MSN	X	Y	CPKpt	SDS/MW	MSN	X	Y	CPKpt	SDS/MW
1008	405	552	-32.3	52,600	1153	921	1158	-13.7	24,700	1246	547	577	-25.3	50,800
1007	1202	848	-7.5	36,500	1154	1564	864	-3.5	35,900	1247	530	576	-26.3	50,900
1008	856	547	-15.0	53,000	1161	637	420	-21.3	68,400	1248	516	572	-27.0	51,200
1009	1284	226	-7.7	123,200	1162	623	397	-21.8	68,800	1250	973	536	-12.7	53,900
1009	996	822	-12.3	37,700	1163	665	397	-20.2	66,700	1251	607	532	-22.4	54,200
1022	1547	403	-4.1	67,900	1168	564	528	-24.4	54,500	1252	665	529	-20.2	54,400
1023	1381	551	-6.4	52,700	1170	552	529	-25.0	54,500	1253	899	766	-14.1	40,200
1024	1525	496	-4.3	57,200	1171	538	534	-25.9	54,800	1254	1311	746	-7.1	41,200
1033	1128	645	-6.7	46,500	1172	545	514	-25.5	55,700	1255	1300	761	-7.5	40,400
1036	1226	274	-4.5	98,300	1174	1099	522	-10.2	55,000	1257	1938	712	0.0	42,900
1039	1761	262	-1.6	103,600	1176	1304	586	-7.5	50,200	1258	1806	718	-1.0	42,600
1040	541	839	-25.7	36,900	1177	1366	539	-6.6	53,700	1259	1727	715	-2.0	42,700
1041	818	910	-15.8	34,000	1178	1608	702	-3.3	43,400	1260	1629	713	-3.0	42,800
1044	1036	485	-11.3	58,300	1179	1485	224	-4.8	124,800	1261	1555	717	-4.0	42,600
1045	1439	407	-5.5	67,300	1180	1459	224	-5.2	124,800	1262	1468	717	-5.0	42,600
1047	1540	250	-4.2	106,200	1181	1431	223	-5.7	125,100	1263	1413	722	-6.0	42,400
1048	1576	535	-3.7	47,100	1182	1407	223	-6.1	125,200	1264	1340	717	-7.0	42,600
1049	1069	411	-10.4	66,700	1183	1383	224	-6.4	124,700	1265	1263	717	-8.0	42,600
1050	949	1040	-13.2	28,900	1184	1454	182	-5.3	164,400	1266	1182	720	-9.0	42,500
1051	426	818	-31.1	37,800	1185	1422	183	-5.8	162,600	1267	1110	717	-10.0	42,600
1052	1583	1385	-3.6	16,900	1186	1394	182	-6.3	164,300	1268	1055	717	-11.0	42,600
1053	779	1092	-16.8	27,000	1189	1171	214	-9.2	131,800	1269	999	717	-12.0	42,600
1054	1613	620	-3.2	48,000	1190	1457	286	-5.2	84,200	1270	959	715	-13.0	42,700
1055	1380	377	-6.5	72,000	1191	686	1114	-19.5	26,200	1271	905	712	-14.0	42,900
1056	284	663	-35.0	45,500	1192	265	863	-35.0	34,700	1272	857	714	-15.0	42,800
1059	1261	746	-8.0	41,200	1193	403	1292	-32.6	20,000	1273	810	705	-16.0	43,300
1060	393	605	-33.3	48,000	1194	344	1275	-35.0	20,600	1274	774	711	-17.0	42,900
1061	1817	645	-0.9	46,600	1195	505	1311	-27.6	19,400	1277	737	708	-18.0	43,100
1062	1245	746	-8.2	41,200	1196	572	1293	-24.1	20,000	1278	702	711	-19.0	42,900
1064	1258	792	-8.1	39,000	1197	639	1502	-21.2	13,000	1279	671	710	-20.0	43,000
1065	705	634	-18.9	33,000	1198	637	1402	-21.3	16,300	1280	645	710	-21.0	43,000
1066	1181	734	-9.0	41,800	1199	614	1407	-22.1	16,200	1281	617	707	-22.0	43,100
1067	529	658	-26.3	45,800	1200	637	1431	-21.3	15,400	1282	595	704	-23.0	43,300
1068	508	696	-27.4	43,700	1201	1095	1394	-10.3	16,600	1283	573	700	-24.0	43,500
1069	1898	604	-0.3	49,100	1202	1719	1545	-2.1	11,600	1284	552	695	-25.0	43,700
1071	873	609	-14.7	48,700	1203	791	668	-16.5	45,200	1285	536	694	-26.0	43,800
1072	1768	1128	-1.5	25,800	1204	964	1021	-12.9	29,700	1286	515	687	-27.0	44,200
1075	836	773	-15.4	39,900	1205	313	195	-35.0	148,700	1287	496	683	-28.0	44,400
1078	1863	861	-0.6	36,000	1208	306	194	-35.0	149,800	1288	467	669	-29.0	45,200
1079	826	566	-15.7	51,600	1209	326	197	-35.0	147,400	1289	447	667	-30.9	45,300
1081	971	833	-12.7	58,500	1210	326	197	-35.0	146,600	1290	427	655	-31.0	45,900
1083	1697	202	-2.3	142,300	1211	394	294	-33.2	91,400	1291	412	655	-32.0	45,900
1085	1157	794	-9.4	38,800	1212	402	294	-32.7	91,200	1292	397	652	-33.0	46,100
1090	620	910	-21.9	34,000	1214	386	294	-33.7	91,400	1293	381	654	-34.0	46,000
1092	1867	597	-0.5	49,500	1215	641	329	-21.2	81,600	1294	365	653	-35.0	46,100
1093	2019	894	-10.0	34,600	1216	660	329	-20.4	81,600					
1094	1546	538	-4.1	53,700	1217	914	266	-13.8	101,800					
1095	1545	477	-4.1	59,100	1218	873	245	-14.7	112,000					
1098	61	935	-35.0	33,000	1219	970	372	-12.7	72,900					
1099	1954	527	-10.0	116,000	1220	1021	298	-11.6	90,100					
1101	588	1048	-22.3	26,600	1221	1392	205	-6.3	139,500					
1102	1050	667	-11.1	45,200	1222	1354	203	-6.8	141,800					
1103	457	797	-29.5	38,800	1223	1362	205	-6.7	139,500					
1105	1884	532	-0.4	54,200	1224	673	540	-19.9	53,600					
1106	1714	648	-2.1	46,300	1225	614	542	-22.1	53,400					
1107	1717	646	-2.1	53,100	1226	603	539	-22.6	53,600					
1108	1976	722	-10.0	42,400	1227	696	623	-19.2	47,800					
1111	547	1066	-25.3	28,000	1228	707	628	-18.9	47,500					
1112	1348	621	-6.9	48,000	1229	475	447	-28.7	62,300					
1115	1385	762	-6.4	40,400	1230	466	1282	-29.0	20,400					
1116	1078	816	-10.6	38,000	1231	759	1461	-17.4	14,400					
1117	975	787	-12.6	39,300	1232	1324	1170	-7.2	24,200					
1118	1022	933	-8.7	33,100	1233	1583	1005	-3.6	30,300					
1119	1022	1076	-11.6	27,600	1234	1865	809	-0.6	38,200					
1120	1005	616	-0.3	48,300	1235	1812	817	-1.0	37,900					
1121	1512	1301	-4.5	19,700	1236	1411	703	-6.0	45,400					
1122	1114	677	-9.9	44,700	1237	1392	682	-6.3	44,500					
1123	1464	452	-5.1	61,700	1238	794	410	-16.4	66,900					
1125	1048	857	-11.1	36,200	1239	769	407	-17.1	67,300					
1126	1122	802	-9.8	36,600	1240	740	406	-17.9	67,500					
1128	1722	892	-2.1	34,700	1241	743	511	-17.8	55,900					
1133	1098	825	-10.2	37,500	1242	713	510	-18.7	56,000					
1139	1830	569	-0.8	51,400	1243	682	509	-19.6	56,100					
1147	764	1182	-17.3	23,800	1244	663	504	-20.3	56,500					
1148	1668	724	-10.0	42,300	1245	565	582	-24.4	50,500					

Table 2. Table of some identified proteins

POP name	Protein name	MSN's	Basis for identification
IDS:3_ALPHA_HDH	3- $\alpha$ -hydroxyisovaleryl-coenzyme A dehydrogenase	137, 159	Pure protein and antibody provided by Dr. T.M. Penning, Department of Pharmacology, School of Medicine, University of Pittsburgh
IDS:ACTIN_BETA	$\beta$ cellular actin, a cytoskeletal protein	38	Homologous position with respect to other mammalian systems
IDS:ACTIN_GAMMA	$\gamma$ cellular actin, a cytoskeletal protein	68	Homologous position with respect to other mammalian systems
IDS:ALBUMIN	Serum albumin, mature form	21, 28, 33	Preference in rat plasma
IDS:APO_A1	Apo A1 plasma lipoprotein, mature form (native)	238, 483	Presence in rat plasma
IDS:CALMODULIN	Calmodulin, lipid, cytosolic calcium-binding protein	123, 649	Homologous position with respect to other mammalian systems
IDS:CATALASE	Catalase (peroxisomal)	54, 61, 106	Presence in purified peroxisomes, similarity in position to mouse catalase
IDS:CPKSPOTS	Spots contributed by the CPK charge standards (muscle proteins)	1257 - 1295	
IDS:CPS	Carbamoyl phosphate synthase	114, 157, 167, 174, 1184, 1185, 1186, 1222	Pure protein provided by Dr. Margaret Marshall, Department of Pharmacology, Medical School, University of Wisconsin - Madison
IDS:CYTOCHROME_B5	Cytochrome b5	87, 477	Pure protein provided by Dr. Andrew Parkinson, Department of Pharmacology, Toxicology and Therapeutics, University of Kansas Medical Center
IDS:FABP_L	Liver fatty acid binding protein	227	Pure protein provided by Dr. Nathan Bass, Department of Medicine, University of California School of Medicine, San Francisco
IDS:HMG-CoA SYNTHASE	Cytosolic HMG-CoA Synthase	133, 144, 235, 413	Antibody provided by Dr. Michael Guarnieri, Merck Sharp & Dohme Research Laboratories, Rahway, NJ
IDS:LAMIN_B	Lamin B, a nuclear protein	415, 734	Homologous position with respect to other mammalian systems
IDS:MITCON:1	Mitcon:1 (F1 ATPase $\beta$ subunit), a mitochondrial inner membrane	17, 49, 71, 340, 1245, 1246, 1247, 1249	Homologous position with respect to other mammalian systems, presence in mitochondria
IDS:MITCON:2	Mitcon:2, a mitochondrial matrix stress	15, 25, 110, 1241, 1242, 1243, 1244	Homologous position with respect to other mammalian systems, presence in mitochondria
IDS:MITCON:3	Mitcon:3, a mitochondrial stress	18, 35, 226, 600, 1238, 1239, 1240	Homologous position with respect to other mammalian systems, presence in mitochondria
IDS:NADPH_P450_RED	NADPH cytochrome P-450 reductase, frequently co-induced with P-450's	175, 251, 812	Pure protein provided by Dr. Andrew Parkinson, Department of Pharmacology, Toxicology and Therapeutics, University of Kansas Medical Center
IDS:PDI	Protein disulphide isomerase 1	168, 1170, 1171, 1172	Sequence information obtained by R.M. Van Frank, Lilly Research Laboratories, Indianapolis
IDS:PLASMA_PROTEINS	Rat plasma proteins observed in liver	21, 28, 33, 44, 72, 102, 115, 197, 238, 246, 248, 257, 293, 332, 347, 364, 365, 419, 432, 433, 468, 518, 552, 605, 623, 665, 687, 725, 738, 780, 805, 833, 926	Plasma electrophoresis studies
IDS:PRO-ALBUMIN	Serum albumin precursor	47, 93	Relative position to mature albumin, presence in micro-Pavica, R.L. et al. BBA (1990) 1029, 115, 125
IDS:PYR-CARBOX	Pyruvate carboxylase	179, 1180, 1181, 1182, 1183	Sequence information obtained by R.M. Van Frank, Lilly Research Laboratories, Indianapolis
IDS:SOD	Superoxide dismutase	56, 132, 1224, 1252	Homologous position with respect to other mammalian systems
IDS:TUBULIN_ALPHA	$\alpha$ tubulin, a cytoskeletal protein	50, 1225, 1226, 1251	Homologous position with respect to other mammalian systems
IDS:TUBULIN_BETA	$\beta$ tubulin, a cytoskeletal protein		

Hb-beta<sub>2</sub>

Computed hemoglobin

Protein

Rabbit r

e 3. Computed pI's of two sets of carbamylated protein standards: Rabbit muscle CPK and human hemoglobin (Hb)

Protein Name	PIR Name	#ASP 3.9	#GLU 4.1	#HIS 6.0	#LYS 10.8	#ARG 12.5	NH2 7.0	Calc pI	Real CPK
Rabbit muscle CPK	KIRBCM	28	27	17	34	18	1	6.84	0.0
		28	27	17	33	18	1	6.67	-1
		28	27	17	32	18	1	6.54	-2
		28	27	17	31	18	1	6.42	-3
		28	27	17	30	18	1	6.31	-4
		28	27	17	29	18	1	6.21	-5
		28	27	17	28	18	1	6.12	-6
		28	27	17	27	18	1	6.03	-7
		28	27	17	26	18	1	5.94	-8
		28	27	17	25	18	1	5.85	-9
		28	27	17	24	18	1	5.76	-10
		28	27	17	23	18	1	5.67	-11
		28	27	17	22	18	1	5.58	-12
		28	27	17	21	18	1	5.48	-13
		28	27	17	20	18	1	5.39	-14
		28	27	17	19	18	1	5.29	-15
		28	27	17	18	18	1	5.20	-16
		28	27	17	17	18	1	5.12	-17
		28	27	17	16	18	1	5.04	-18
		28	27	17	15	18	1	4.96	-19
		28	27	17	14	18	1	4.89	-20
		28	27	17	13	18	1	4.83	-21
		28	27	17	12	18	1	4.77	-22
		28	27	17	11	18	1	4.71	-23
		28	27	17	10	18	1	4.66	-24
		28	27	17	9	18	1	4.61	-25
		28	27	17	8	18	1	4.56	-26
		28	27	17	7	18	1	4.52	-27
		28	27	17	6	18	1	4.48	-28
		28	27	17	5	18	1	4.44	-29
		28	27	17	4	18	1	4.40	-30
		28	27	17	3	18	1	4.36	-31
		28	27	17	2	18	1	4.32	-32
		28	27	17	1	18	1	4.29	-33
		28	27	17	0	18	1	4.25	-34
		28	27	17	0	18	0	4.22	-35
Hb-beta, human	HBHU	7	8	9	11	3	1	7.18	
		7	8	9	10	3	1	6.79	
		7	8	9	9	3	1	6.53	-1.8
		7	8	9	8	3	1	6.32	-3.2
		7	8	9	7	3	1	6.13	-5.3
		7	8	9	6	3	1	5.96	-7.2
		7	8	9	5	3	1	5.78	-10.0
		7	8	9	4	3	1	5.59	-12.3
		7	8	9	3	3	1	5.37	-15.5
		7	8	9	2	3	1	5.14	-18.0
		7	8	9	1	3	1	4.91	-21.0
		7	8	9	0	3	1	4.71	-25.5
		7	8	9	0	3	0	4.54	-27.2

Table 4. Computed pI's of some known proteins related to measured CPK pI's

Protein Name	PIR Name	#ASP 3.9	#GLU 4.1	#HIS 6.0	#LYS 10.8	#ARG 12.5	Calc pI	Real CPK
0 Creatine phospho kinase (CPK), rabbit muscle	KIRBCM	28	27	17	34	18	6.84	0.0
1 Fatty acid-binding protein, rat hepatic	FZRTL	5	13	2	16	2	7.83	-3.0
2 b2-microglobulin, human	MGHUB2	7	8	4	8	5	6.09	-5.0
3 Carbamoyl-phosphate synthase, rat	SYRTCA	72	96	28	95	56	5.97	-5.5
4 Proalbumin (serum albumin precursor), rat	ABRTS	32	57	15	53	27	5.96	-6.2
5 Serum albumin, rat	ABRTS	32	57	15	53	24	5.71	-9.0
6 Superoxide dismutase (Cu-Zn, SOD), rat	A26810	8	11	10	9	4	5.91	-9.2
7 Phospholipase C, phosphoinositide-specific (?), rat	A28807	34	42	9	49	21	5.92	-9.2
8 Albumin, human	ABHUS	36	61	16	60	24	5.70	-11.9
9 Apo A-I lipoprotein, rat	A24700	18	24	6	23	12	5.32	-13.7
10 proApo A-I lipoprotein, human	LPHUA1	16	30	6	21	17	5.35	-14.3
11 NADPH cytochrome P-450 reductase, rat	RDRT04	41	60	21	38	36	5.07	-15.6
12 Retinol binding protein, human	VAHU	18	10	2	10	14	5.04	-16.9
13 Actin beta, rat	ATRTC	23	26	9	19	18	5.06	-17.2
14 Actin gamma, rat	ATRTC	20	29	9	19	18	5.07	-16.6
15 Apo A-I lipoprotein, human	LPHUA1	16	30	5	21	16	5.10	-17.5
16 Apo A-IV lipoprotein, human	LPHUA4	20	49	8	28	24	4.88	-19.7
17 Tubulin alpha, rat	UBRTA	27	37	13	19	21	4.66	-19.8
18 F1ATPase beta, bovine	PWB08	25	36	9	22	22	4.80	-21.0
19 Tubulin beta, pig	UBPGB	26	36	10	15	22	4.49	-22.5
20 Protein disulphide isomerase (PDI), rat hepatic	ISRTSS	43	51	11	51	9	4.07	-25.0
21 Cytochrome b5, rat	CBRT5	10	15	6	10	4	4.59	-26.0
22 Apo C-II lipoprotein, human	LPHUC2	4	7	0	6	1	4.44	-30.5
Amino acid pI assumed in calculation:		3.9	4.1	6.0	10.8	12.5		

Wirth  
Lao  
Fujimoto  
C. Bisgaard  
D. Olson  
History of Exp.  
ogenesis.  
Cancer In  
Institutes  
ada,

ents  
roduction  
aterials and  
Materials.  
Cells.  
Metabolic  
nine and  
Sample p:  
Subcellul:  
2-D PAG  
Computer  
retograms  
ults  
[<sup>35</sup>S]Methi  
Whole ce  
2Subcellul.  
[<sup>32</sup>P]Ortho  
discussion  
ferences  
endum 1:  
endum 2:  
eins

vidence: Dr. P  
National Car  
USA

ons: 2-D PA  
HLE, hum.  
weight; NE  
Nonidet P  
; RLE, rat  
tagewellch.



N. Leigh Anderson<sup>1</sup>  
 Ricardo Esquer-Blasco<sup>2</sup>  
 Jean-Paul Hofmann<sup>1</sup>  
 Lydie Meheurs<sup>1</sup>  
 Jos Raymakers<sup>1</sup>  
 Sandra Steiner<sup>3</sup>  
 Frank Witzmann<sup>4</sup>  
 Norman G. Anderson<sup>1</sup>

<sup>1</sup>Large Scale Biology Corporation,  
 Rockville, MD  
<sup>2</sup>Innogenetics NV, Ghent  
<sup>3</sup>Sandoz Pharma Ltd, Drug Safety  
 Assessment, Toxicology, Basel  
<sup>4</sup>Molecular Anatomy Laboratory,  
 Indiana University Purdue  
 University Columbus, Columbus, IN

## An updated two-dimensional gel database of rat liver proteins useful in gene regulation and drug effect studies

We have improved upon the reference two-dimensional (2-D) electrophoretic map of rat liver proteins originally published in 1991 (N. L. Anderson *et al.*, *Electrophoresis* 1991, 12, 907-930). A total of 53 proteins (102 spots) are now identified, many by microsequencing. In most cases, spots cut from wet, Coomassie Blue stained 2-D gels were submitted to internal tryptic digestion [2], and individual peptides, separated by high-performance liquid chromatography (HPLC), were sequenced using a Perkin-Elmer 477A sequencer. Additional spots were identified using specific antibodies.

Figure 1 shows the current annotated 2-D map of F344 rat liver, analyzed using the Iso-DALT system (20 × 25 cm gels) and BDH 4-8 carrier ampholytes. Both the map itself and the master spot number system remain the same as shown in the original publication. Table 1 lists the important features of each identification shown, including the gel position,  $pI$ , and  $M_r$ , for the most abundant or most basic form of each protein. Using this extended base of identified spots, a series of four improved calibration functions has been derived for the  $pI$  and SDS- $M_r$  axes (the first two of which are shown in Fig. 2A and B). Both forward and reverse functions are derived, so that one can compute the physical properties of a spot with a given gel location, or inversely compute the gel position expected for a protein having given physical properties:

$$Y_{\text{RATLIVER}} = f_{\text{M}}(\text{RATLIVER}, M_r(\text{SEQUENCE-DERIVED})) \quad (1)$$

$$X_{\text{RATLIVER}} = f_{\text{pI}}(\text{RATLIVER}, pI(\text{SEQUENCE-DERIVED})) \quad (2)$$

$$M_r(\text{GEL-DERIVED}) = f_{\text{RATLIVER}}^{-1}(M_r(Y_{\text{RATLIVER}})) \quad (3)$$

$$pI(\text{GEL-DERIVED}) = f_{\text{RATLIVER}}^{-1}(pI(X_{\text{RATLIVER}})) \quad (4)$$

A spreadsheet program (in Microsoft Excel) was developed to facilitate flexible computation of  $pI$ 's from amino acid sequence data, and the results were entered into a relational database (Microsoft Access). A table of spot positions and sequence-derived  $pI$ 's and  $M_r$ 's was fitted with a large series of analytic equations using Tablecurve (Jandel Scientific), and the four conversion Eqs. (1)-(4), relating computed  $pI$  and gel  $X$  coordinate, or computed molecular weight and gel  $Y$  coordinate, were selected, based on criteria of simplicity, goodness of fit and favorable asymptotic behavior. Table 2 lists the equations and coefficients. Application of Eqs. (3) and (4) to a spot's  $X$  and  $Y$  coordinates, given in [1], produce improved  $M_r$  estimates, and allow computation of  $pI$

directly in pH units, instead of in terms of positions relative to creatine phosphokinase (CPK) charge standards. The inverse Eqs. (1) and (2) were used to compute the gel positions of a series of  $pI$  and  $M_r$  tick marks. These tick marks were plotted with SigmaPlot (Jandel), together with fiducial marks locating several prominent spots, and the resulting graphic was aligned over the synthetic gel image (computed by Kepler from the master gel pattern) using Freehand (Lotus Development). Maps were printed as Postscript output from Freehand, either in black and white (as shown here) or in color, where label color indicates subcellular location (available from the first author upon request). We have also used the rat liver 2-D pattern as presented here to calibrate the patterns of other samples. Using mixtures of rat liver and mouse liver samples, for example, we made composite 2-D patterns that allow use of the rat pattern to standardize both axes of the mouse pattern. This was accomplished by deriving transformations relating the rat and mouse  $X$ , and separately the rat and mouse  $Y$ , axes (Table 2, lower half; Fig. 2C and D) based on a series of spots that coelectrophorese in these closely related species. These functions were then applied to derive equations relating the mouse liver  $X$  and  $Y$  to  $pI$  and SDS- $M_r$  (Eqs. 5 and 6 below). The resulting standardized 2-D pattern for B6C3F1 mouse liver is shown in Fig. 3.

$$M_r(\text{MOUSELIVER}) = f_{\text{RATLIVER}}^{-1}(M_r(\text{MOUSELIVER} - \text{RATLIVER} Y (Y_{\text{MOUSELIVER}}))) \quad (5)$$

$$pI(\text{MOUSELIVER}) = f_{\text{RATLIVER}}^{-1}(pI(\text{MOUSELIVER} - \text{RATLIVER} X (X_{\text{MOUSELIVER}}))) \quad (6)$$

A slightly more complex approach can be used to standardize samples that have few or no spots co-electrophoresing with rat liver proteins. In this case, a 2-D gel is prepared with a mixture of the two samples, and four functions (forward and backward, each for  $X$  and  $Y$ ) are derived relating each sample's own master pattern to the composite. The required functions are then applied in a nested fashion to yield the desired result (using rat plasma as an example):

$$M_r(\text{RATPLASMA}) = f_{\text{RATLIVER}}^{-1}(M_r(\text{RATPLASMA} - \text{LIVER} Y (\text{RATPLASMA} - \text{LIVER} Y (Y_{\text{RATPLASMA}})))) \quad (7)$$

Correspondence: Dr. Leigh Anderson, Large Scale Biology Corporation, 9620 Medical Center Drive, Rockville, MD 20850-3338 USA (Tel: +301-424-5989; Fax: +301-762-4892; email: leigh@lsbc.com)

Keywords: Two-dimensional polyacrylamide gel electrophoresis / Liver / Map / Identification / Calibration

F344 RAT LIVER 2-D PROTEIN PATTERN

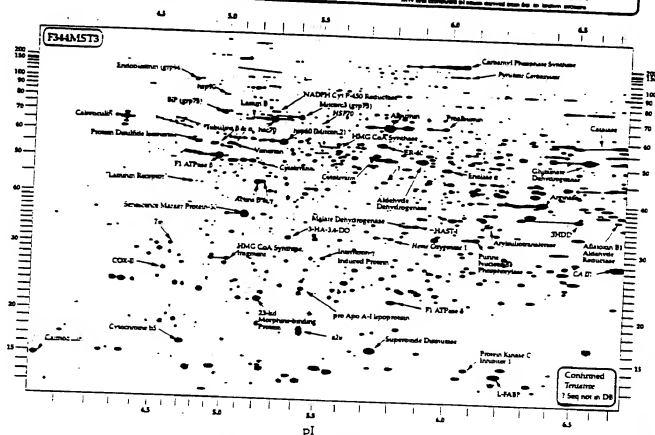


Figure 1. Master 2-D gel pattern of Fischer 344 rat liver proteins, annotated with 53 protein identifications and computed *pI* and *M<sub>r</sub>* axes. Tentative identifications are in italic type.

**Table 1. Proteins identified in the 2-D pattern of F344 rat liver**

<sup>15</sup> N	Protein ID <sup>a</sup>	Protein name	Identification comments	GeI X <sup>b</sup>	Experimental pI <sup>b</sup>	GeI Y <sup>c</sup>	Experimental M <sub>r</sub> <sup>c</sup>
126	HADO_HUMAN <sup>a</sup>	3-HA-3,4-DO: 3-hydroxy-anthranilate 3,4-dioxygenase	Internal sequence	871.95	5.36	921.35	30 207
137, 159, 228, 258	DIDH_RAT	3HDD: 3-hydroxysteroid dihydrodiol reductase	Ab (T.M. Penning) and pure protein	1857.52	6.51	822.52	34 406
173	MUP_RAT	$\alpha$ u globulin	Presence in liver microsome lumen.	919.16	5.43	1313.81	19 549
38	ACTB_HUMAN	Actin $\beta$	Abundance in kidney; pI, M <sub>r</sub>	763.40	5.19	693.64	41 586
68	ACTG_HUMAN	Actin $\gamma$	Analogy with other mammalian patterns (e.g. human) through coelectrophoresis	779.42	5.21	692.26	41 593
693	AFAR_RAT	Alfalozin B1 aldehyde	Analogy with other mammalian patterns (e.g. human) through coelectrophoresis	1993.32	6.72	818.60	34 677
28, 21, 33	ALBU_RAT	Albumin	Internal sequence	1262.81	5.86	445.64	66 354
43	DHAM_RAT	Aldehyde dehydrogenase	Coelectrophoresis with principal plasma protein	1317.72	5.91	589.03	49 602
96	ARGI_RAT	Arginase	N-Terminal sequence and AAA	1730.72	6.34	756.02	37 819
117	SUAR_RAT	Arginylsulfotransferase	Internal sequence	1547.96	6.14	849.08	33 186
163, 1161, 1162, 20	GR7B_RAT	BIP (GRP-78)	Ab (F. Witzmann)	665.33	5.01	397.39	74 564
1185	CAH3_RAT	CA-III	Uncertain; by comparison with mouse	1996.60	6.72	1017.02	26 887
122	CALM_HUMAN	Calmodulin	Analogy with human cellular patterns through coelectrophoresis	23.05	4.03	1433.25	17 419
3, 201, 48, 39, 22, 24	CRTC_RAT	Calreticulin	Ab (Lance Pohl)	310.59	4.34	433.80	68 206

Table 1. continued

MSN <sup>a</sup>	Protein ID(s)	Protein name	Identification comments	Gel X <sup>b</sup>	Experimental pI <sup>c</sup>	Gel Y <sup>b</sup>	Experimental M <sup>d</sup>
1184, 1186, 114, 174, 118, 5, 167, 157	CPSM_RAT	Carbamoyl phosphate synthase	2-D of pure protein; confirmed by N-terminal sequence and AAA	143.56	6.05	181.64	160 640
54, 61	CATA_RAT	Catalase	Internal sequence	2000.81	6.73	499.64	58 968
136	COX2_RAT	COX-II	Ab (J. W. Taanman), confirmed by internal sequence	452.57	4.61	1062.67	25 504
87	CYB5_RAT	Cytochrome B5	2-D of pure protein; Ab; confirmed by AAA	515.68	4.73	1370.55	18 493
41	CK-RAT <sup>1</sup>	Cytokeratin	Location in cytoskeletal fraction	1165.12	5.75	569.09	51 448
29	CK-RAT <sup>2</sup>	Cytokeratin	Location in cytoskeletal fraction	743.11	5.15	605.23	48 187
5, 11	ENPL-RAT <sup>1</sup>	Endoplasmic	Ab (F. Witzmann)	567.70	4.83	263.37	112 194
60	ENO4_RAT	Enolase A	Internal sequence and AAA	1399.78	6.00	623.54	46 674
27	ER60_RAT	ER-60	N-Terminal sequence (R. M. Van Frank)	1184.20	5.77	523.51	56 169
17	ATPB_RAT	F1 ATPase B	N-Terminal sequence and AAA	629.06	4.95	588.83	49 620
196	ATP7_RAT	F1 ATPase A	Internal sequence	1227.24	5.82	1184.65	22 310
79	F16P_RAT	Fructose-1,6-bis-phosphatase	Uncertain; by comparison with 1D in Garrison and Wager (JBC 257:13135-13143)	924.54	5.44	737.77	38 858
62, 78	DHE3_RAT	Glutamate dehydrogenase	N-Terminal sequence and internal sequence	1887.39	6.55	566.92	51 655
125	HAST-RAT <sup>1</sup>	HAST-1: N-hydroxyarylamine sulfoxidase	Internal sequence	1297.94	5.89	861.55	32 638
307	HO1_RAT	Heme oxygenase 1	Uncertain; available data from internal sequence	1219.39	5.81	915.71	30 423
413, 1250, 933	HMCS_RAT	HMG CoA synthase, cytosolic	Ab (J. Gernemshausen)	1033.48	5.59	538.13	54 571
133, 144, 235	HMCS_RAT	HMG CoA synthase, mitochondrial (frag)	Ab (J. Gernemshausen), N-terminal sequence (Steiner/Lotzspeich)	666.40	5.02	1019.42	26 811
8, 23, 1307	H57C_RAT	HSC-70	Positional homology; (with human, etc.) through coelectrophoresis	811.87	5.27	425.76	69 521
15, 25, 110	P60_RAT	HSP-60	Ab (F. Witzmann), confirmed by N-terminal sequence and AAA	845.09	5.32	520.03	56 561
971	H570-RAT <sup>1</sup>	HSP-70	Ab (F. Witzmann)	976.11	5.51	437.14	67 674
1216, 1215, 90	H590-RAT <sup>1</sup>	HSP-90	Ab (F. Witzmann)	659.86	5.00	329	90 107
256	ING1-HUMAN	Interferon- $\gamma$ induced protein	Internal sequence	993.85	5.54	1006.04	27 237
415, 734	LAMB-RAT <sup>1</sup>	Lamin B	Positional homology with human through coelectrophoresis, nuclear location	737.10	5.14	425.19	69 615
80	LAMR-RAT <sup>1</sup>	"Laminin receptor"	Internal sequence	534.02	4.77	697.62	41 327
227	FABL_RAT	L-FABP (liver fatty acid binding protein)	Ab (N. M. Bass)	1586.09	6.18	1483.43	16 622
134	MDHC_MOUSE	Malate dehydrogenase	Internal sequence	1270.85	5.86	861.96	32 620
18, 35, 226	GR75-RAT <sup>1</sup>	Mitochondrion: p75	Positional homology with human through coelectrophoresis	905.67	5.41	413.67	71 589
175, 251	NCPR_RAT	NADPH P450 reductase	2-D of pure protein	824.69	5.29	393.21	75 366
1168, 1170, 1171	PDI_RAT	PDI: Protein disulfide isomerase	N-Terminal sequence (R. M. van Frank), Ab	564.30	4.83	528.47	55 618
47, 93	ALBU_RAT	Pro-Albumin	Microsomal lumen location, pI, M <sub>r</sub> relative to albumin	1391.03	5.99	446.68	66 195
236	APA1_RAT	Pro-APO A-I lipoprotein	Coelectrophoresis with plasma protein	920.41	5.43	1177.51	23 467
320	IPK1_BOVIN	Protein kinase C inhibitor	Internal sequence; homology with bovine protein	1480.01	6.08	1458.81	17 007
152	PNFH_MOUSE	Purine nucleoside phosphorylase	Internal sequence	1507.19	6.10	911.16	30 599
1179, 1180, 1181, 1182, 1183	PYVC-RAT <sup>1</sup>	Pyruvate carboxylase	Tentative; 2-D of pure protein (J. G. Henzel, JBC, 1979); reported in <i>Biochim Biophys. Acta</i> 1022, 115-125	1485.10	6.08	223.52	131 589
55, 103	SM30_RAT	SMP-30: Serine/threonine marker protein-30	Internal sequence	721.71	5.11	830.10	34 051
135	SODC_RAT	Superoxide dismutase	AAA; confirmed by internal sequence (R. M. Van Frank)	1161.24	5.74	1388.68	18 173
172	TPM-RAT <sup>1</sup>	Tm: tropomyosin	Location in cytoskeleton, 2-D position relative to human, Ab	476.24	4.66	957.86	28 865
277, 56	TBA1_RAT	Tubulin $\alpha$	Positional homology with human through coelectrophoresis, cytoskeletal location	688.22	5.06	537.67	54 620
50, 1225	TBB1_RAT	Tubulin $\beta$	Positional homology with human through coelectrophoresis, cytoskeletal location	621.29	4.93	535.48	54 855
1224	VIME_RAT	Vimentin	Positional homology with human through coelectrophoresis, cytoskeletal location	673.00	5.03	539.50	54 426

Table 1. continued

MSN <sup>(1)</sup>	Protein ID(s)	Protein name	Identification comments	Gel X <sup>(2)</sup>	Experimental pI <sup>(3)</sup>	Gel Y <sup>(4)</sup>	Experimental M <sub>r</sub> <sup>(5)</sup>
5	Unknown	? not in sequence databases	Internal sequence	1191.28	5.78	680.42	47 469
	BBPL_RAT	23 kDa morphine-binding protein	Internal sequence	773.31	5.20	1182.41	22 363

<sup>(1)</sup> Master spot number (MSN) from [1]

<sup>(2)</sup> PROT identifier

<sup>(3)</sup> coordinates of the most basic or most abundant assigned spot on the F344 master gel pattern

<sup>(4)</sup> and <sup>(5)</sup> of the most basic or most abundant assigned spot, derived from the calibration functions included here

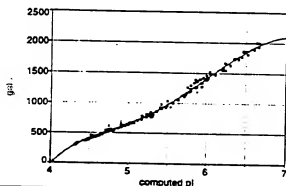
<sup>(6)</sup> PROT style proposed identifier

Abbreviations: AAA, amino acid analysis; Ab, antibody

## 2. Equations and coefficients

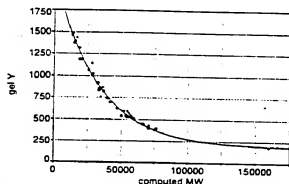
eq	Equation (f)	r <sup>2</sup>	a	b	c	d	e
Y = f(computed M <sub>r</sub> )	$y = a - b \exp(-x/c)$	0.988181021	178.74803	1967.7892	33363.958		
X = f(computed pI)	$y = a - bx - cx/\ln x - dx + e/x^{1.5}$	0.99247216	-868566.5	-904497.94	3856926.1	182768.44	-27154534
cal M <sub>r</sub> = f(rat gel Y)	$y = a - bxc$	0.9960177	-8464.5809	19095881	-0.9086255		
cal pI = f(rat gel X)	$y = a + bx + cx^2 + dx^3 \ln x + ex^3$	0.99176499	4.044686	-0.00114238	0.0000323	-0.000000455	0.00000000176
gel Y = f(rat gel Y)	$y = a + bx + cx^{1.5} + dx^{2.5} \ln x + ex/\ln x$	0.99951069	11861.44	678.91666	-0.78964914	1567.5639	-6953.9592
gel X = f(rat gel X)	$y = a + bx^2 \ln x + cx^{2.5} + dx^3$	0.99926349	58.935923	0.00091353	-0.000213688	0.00000159	
Y = f(mouse gel Y)	$y = a + bx^2 \ln x + cx^{2.5} + dx^3$	0.99950032	69.740526	0.00050772	-0.000130392	0.00000116	
X = f(mouse gel X)	$y = a + bx + cx^2 \ln x + dx^{2.5} + ex^3$	0.9992832	-198.07189	2.0899063	-0.000671191	0.000145189	-0.000000986

$$y = a + bx + cx/\ln x + dx + e/x^{1.5}$$



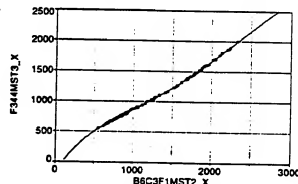
B

$$y = a + b \exp(-x/c)$$



C

$$y = a + bx + cx^2 \ln x + dx^3 + ex^3$$



D

$$y = a + bx^2 \ln x + cx^2 + dx^3$$

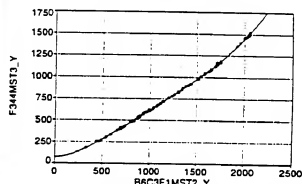
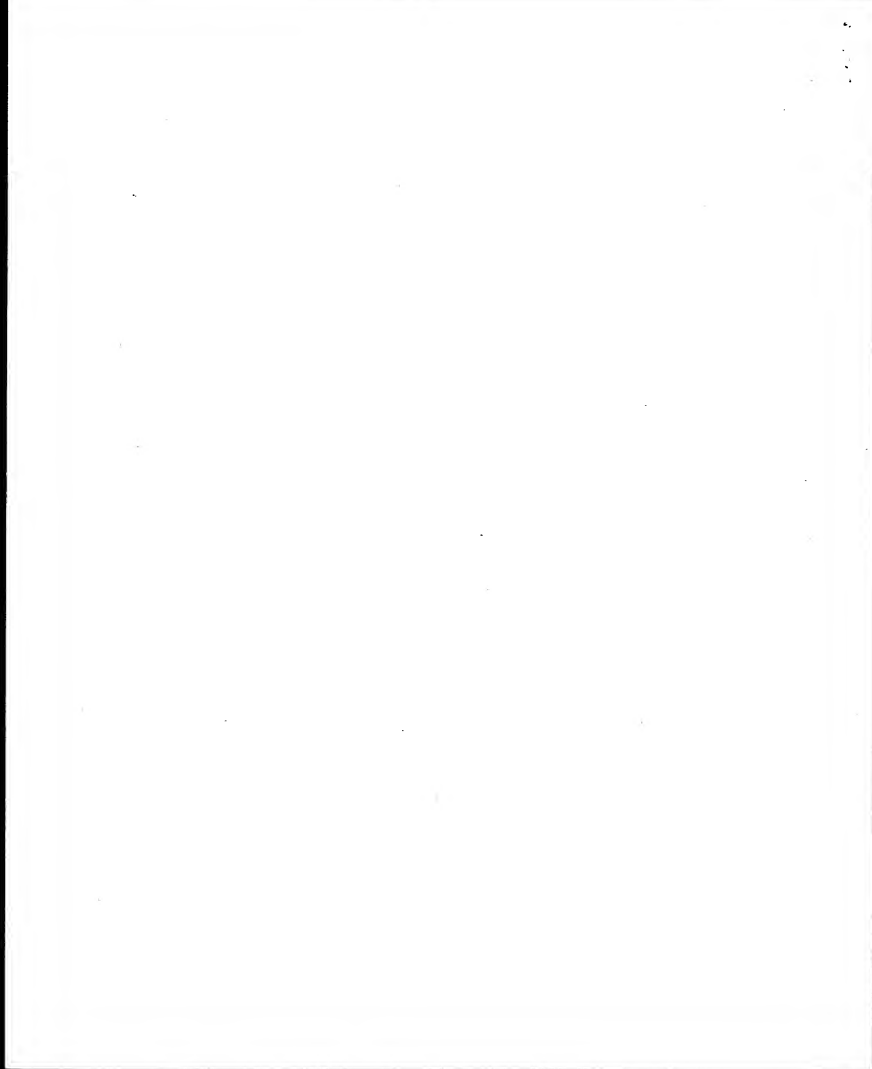


Figure 2. Plots showing fits of selected equations (continuous curves) to data on identified proteins (square symbols). (A) pI computed from sequence data versus gel X position for identified spots in F344 rat liver; (B) M<sub>r</sub> computed from sequence data versus gel Y position for identified spots in F344 rat liver; (C) gel X position for spots in B6C3F1 mouse liver versus X position in F344 rat liver, for coelectrophoresing spots; (D) gel Y position for spots in B6C3F1 mouse liver versus Y position in F344 rat liver, for coelectrophoresing spots. In each case, inverse equations were also computed (Table 2).





# Progress with Proteome Projects: Why all Proteins Expressed by a Genome Should be Identified and How To Do It

MARC R. WILKINS<sup>1</sup>, JEAN-CHARLES SANCHEZ<sup>1</sup>, ANDREW A. GOOLEY<sup>2</sup>,  
RON D. APPEL<sup>3</sup>, IAN HUMPHERY-SMITH<sup>3</sup>, DENIS F. HOCHSTRASSER<sup>1</sup>  
AND KEITH L. WILLIAMS<sup>1,\*</sup>

<sup>1</sup>Macquarie University Centre for Analytical Biotechnology, Macquarie University, Sydney, NSW 2109, Australia; <sup>2</sup>Department of Microbiology, University of Sydney, NSW 2006, Australia and <sup>3</sup>Central Clinical Chemistry Laboratory, and Medical Computing Centre of the University of Geneva, CH 1211 Geneva 14, Switzerland

## Introduction

The advent of large genome sequencing projects has changed the scale of biology. Over a relatively short period of time, we have witnessed the elucidation of the complete nucleotide sequence for bacteriophage  $\lambda$  (Sanger *et al.*, 1982), the nucleotide sequence of a eukaryotic chromosome (Oliver *et al.*, 1992), and in the near future will see the definition of all open reading frames of some simple organisms, including *Mycoplasma pneumoniae*, *Escherichia coli*, *Saccharomyces cerevisiae*, *Cuennihabditus elegans* and *Arabidopsis thaliana*. Nevertheless, genome sequencing projects are not an end in themselves. In fact, they only represent a starting point to understanding the function of an organism. A great challenge that biologists now face is how the co-expression of thousands of genes can best be examined under physiological and pathophysiological conditions, and how these patterns of expression define an organism.

There are two approaches that can be used to examine gene expression on a large scale. One uses nucleic acid-based technology, the other protein-based technology. The most promising nucleic-acid based technology is differential display of mRNA (Liang and Purdee, 1992; Bauer *et al.*, 1993), which uses polymerase chain reaction with arbitrary primers to generate thousands of cDNA species, each which correspond to an expressed gene or part of a gene. However, it is currently unclear if this technique can be developed to reliably assay the expression of thousands of genes or

\* Corresponding Author

identify all cDNA species, and the approach does not easily allow a systematic screening. Analysis of gene expression by the study of proteins present in a cell or tissue presents a favorable alternative. This can be achieved by use of two-dimensional (2-D) gel electrophoresis, quantitative computer image analysis, and protein identification techniques to create 'reference maps' of all detectable proteins. Such reference maps establish patterns of normal and abnormal gene expression in the organism, and allow the examination of some post-translational protein modifications which are functionally important for many proteins. It is possible to screen proteins systematically from reference maps to establish their identities.

To define protein-based gene expression analysis, the concept of the 'proteome' was recently proposed (Wilkins *et al.*, 1995; Wasinger *et al.*, 1995). A proteome is the entire PROTein complement expressed by a genome, or by a cell or tissue type. The concept of the proteome has some differences from that of the genome, as while there is only one definitive genome of an organism, the proteome is an entity which can change under different conditions, and can be dissimilar in different tissues of a single organism. A proteome nevertheless remains a direct product of a genome. Interestingly, the number of proteins in a proteome can exceed the number of genes present, as protein products expressed by alternative gene splicing or with different post-translational modifications are observed as separate molecules on a 2-D gel. As an extrapolation of the concept of the 'genome project', a 'proteome project' is research which seeks to identify and characterize the proteins present in a cell or tissue and define their patterns of expression.

Proteome projects present challenges of a similar magnitude to that of genome projects. Technically, the 2-D gel electrophoresis must be reproducible and of high resolution, allowing the separation and detection of the thousands of proteins in a cell. Low copy number proteins should be detectable. There should be computer gel image analysis systems that can qualitatively and quantitatively catalog the electrophoretically separated proteins, to form reference maps. A range of rapid and reliable techniques must be available for the identification and characterisation of proteins. As a consequence of a proteome project, protein databases must be assembled that contain reference information about proteins; such databases must be linked to genomic databases and protein reference maps. Databases should be widely accessible and easy to use.

Recently, there have been many changes in the techniques and resources available for the analysis of proteomes. It is the aim of this chapter to discuss the status of the areas outlined above, and to review briefly the progress of some current proteome projects.

## Two-dimensional electrophoresis of proteomes

Two dimensional (2-D) gel electrophoresis involves the separation of proteins by their isoelectric point in the first dimension, then separation according to molecular weight by sodium dodecyl sulfate electrophoresis in the second dimension. Since first described (Klose, 1975; O'Farrell, 1975; Scheele, 1975), it has become the method of choice for the separation of complex mixtures of proteins, albeit with many modifications to the original techniques. 2-D electrophoresis forms the basis of proteome projects through separating proteins by their size and charge (Hochstrasser *et al.*,

Fig.  
11a  
was  
top  
the  
of 1

19  
pro  
sin

2-D  
A;  
ph  
en



## HEPG2 2D-PAGE MAP

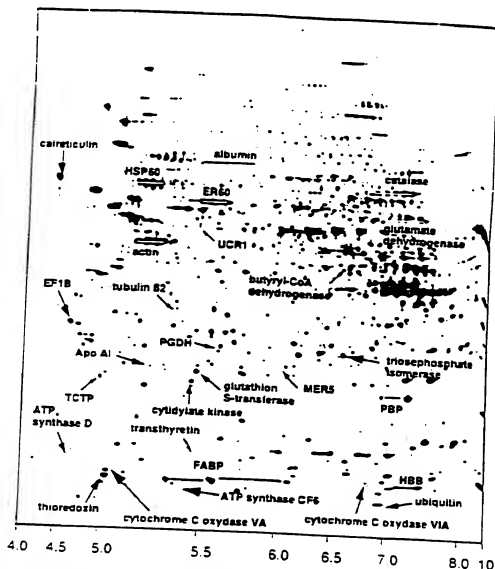


Figure 1. Two-dimensional gel electrophoresis map of a human hepatoblastoma-derived cell line, illustrating the very high resolution of the technique. The first dimensional separation (right to left of figure) was achieved using immobilised pH gradient electrophoresis of 4.0 to 10.0 units. The second dimension (top to bottom of figure) was SDS-PAGE using a 11%–14% acrylamide gradient, allowing separation in the molecular weight range 10–250 kDa. Proteins were visualised by silver staining. Arrows show proteins of known identity.

1992; Celis *et al.*, 1993; Garrels and Franza, 1989; VanBogelen *et al.*, 1992). Current protocols can resolve two to three thousand proteins from a complex sample on a single gel (Figure 1).

#### 2-D GEL RESOLUTION AND REPRODUCIBILITY

A primary challenge of separating complex mixtures of proteins by 2-D gel electrophoresis has been to achieve high resolution and reproducibility. High resolution ensures that a maximum of protein species are separated, and high reproducibility is

vital to allow comparison of gels from day to day and between research sites. These factors can be difficult to achieve.

Carrier ampholytes are a common means of isoelectric focusing for the first dimension of 2-D electrophoresis. Gels are usually focused to equilibrium to separate proteins in the pI range 4 to 8, and run in a non-equilibrium mode (NEPHGE) to separate proteins of higher pI (7 to 11.5) (O'Farrell, 1975; O'Farrell, Goodman and O'Farrell, 1977). Unfortunately, the use of carrier ampholytes in the isoelectric focusing procedure is susceptible to 'cathode drift', whereby pH gradients established by pre-focusing of ampholytes slowly change with time (Righetti and Drysdale, 1973). Carrier ampholyte pH gradients are also distorted by high salt concentration of samples (Bjellqvist *et al.*, 1982), and by high protein load (O'Farrell, 1975). A further limitation is that isoelectric focusing gels, which are cast and subject to electrophoresis in narrow glass tubes, need to be extruded by mechanical means before application to the second dimension - a procedure that potentially distorts the gel. Nevertheless, many of the above shortcomings can be avoided by loading small amounts of  $^{14}\text{C}$  or  $^{35}\text{S}$  radiolabelled samples (Garrels, 1989; Neidhardt *et al.*, 1989; Vandekerckhove *et al.*, 1990). High sensitivity detection is then achieved through use of fluorography or phosphorimaging plates (Bonner and Laskey, 1974; Johnston, Pickett and Barker, 1990; Patterson and Lutter, 1993). However, this approach is only practicable for organisms or tissues that can be radiolabelled.

An alternative technique, which is becoming the method of choice for the first dimension separation of proteins, involves isoelectric focusing in immobilized pH gradient (IPG) gels (Bjellqvist *et al.*, 1982; Gorg, Postel and Gunther, 1988; Righetti, 1990). Immobilized pH gradients are formed by the covalent coupling of the pH gradient into an acrylamide matrix, creating a gradient that is completely stable with time. IPG gels are usually poured onto a stiff backing film, which is mechanically strong and provides easy gel handling (Ostergren, Eriksson and Bjellqvist, 1988). The major advantages of IPG separations are that they do not suffer from cathodic drift, they allow focusing of basic and very acidic proteins to equilibrium, pH gradients can be precisely tailored (linear, stepwise, sigmoidal), and that separations over a very narrow pH range are possible (0.05 pH units per cm) (Righetti, 1990; Bjellqvist *et al.*, 1982, 1993a; Sinha *et al.*, 1990; Gorg *et al.*, 1988; Gelfi *et al.*, 1987; Gunther *et al.*, 1988). However, it is not currently possible to use IPG gels to separate very basic proteins of isoelectric point greater than 10, although this is under development. Narrow pH range separations are useful to address problems of protein co-migration in complex samples, allowing 'zooming in' on regions of a gel (Figure 2). IPG gel strips are now commercially available, which begin to address the problems of intra- and inter-lab isoelectric focusing reproducibility.

There are two means of electrophoresis for the second dimension separation of proteins: vertical slab gels and horizontal ultrathin gels (Gorg, Postel, and Gunther, 1988). Both are usually SDS-containing gradient gels of approximately 11% to 15% acrylamide, which separate proteins in the molecular mass range of 10 - 150 kD. A stacking gel is not usually used with slab gels, but is necessary when using horizontal gel setups (Gorg, Postel and Gunther, 1988). Comparisons have shown that there is little or no difference in the reproducibility of electrophoresis using either approach (Corbett *et al.*, 1994a), but commercially available vertical or horizontal precast gels will provide greater reproducibility for occasional users. For slab gel electrophoresis,

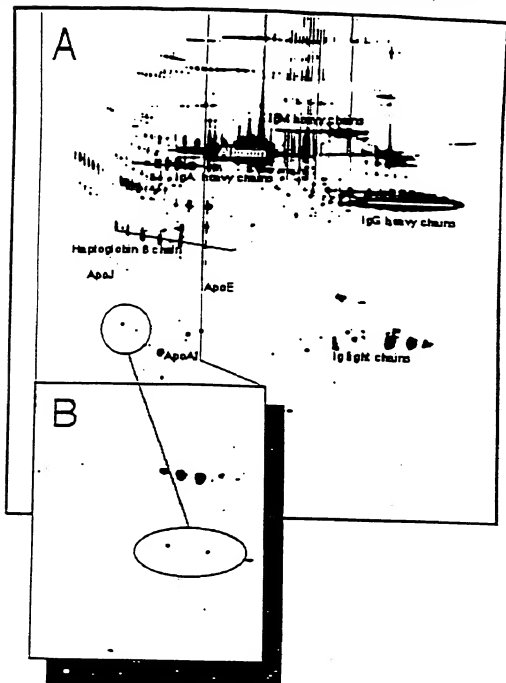


Figure 2. Two-dimensional gel electrophoresis allows 'zooming in' on areas of interest. Rings highlight 2 proteins common to each gel. (A) Wide pI range two dimensional electrophoresis map of human plasma proteins. First dimension separation was achieved using an immobilised pH gradient of 3.5 to 10.0 units. The second dimension was SDS-PAGE. Actual gel size was 16cm x 20cm, and proteins were visualised with silver staining. (B) Narrow pI range electrophoresis was used to 'zoom in' on a small region of the plasma map. The first dimension used a narrow range immobilised pH gradient of 4.2 to 5.2 units, and second dimension was SDS-PAGE. Micropreparative loading was used, and the gel blotted to PVDF. Proteins were visualised with amido black. Actual blot size was 16cm x 20cm.

the use of piperazine diacrylyl as a gel crosslinker and the addition of thiosulfate in the catalyst system has been shown to give better resolution and higher sensitivity detection (Hochstrasser and Merril, 1988; Hochstrasser, Patchornik and Merril, 1988).

Notwithstanding the advances described above, there is an increasing demand to improve the reproducibility of 2-D electrophoresis to facilitate database construction and proteome studies. Harrington *et al.* (1993) explain that if a gel resolves 4000 protein spots, and there is 99.5% spot matching from gel to gel, this will produce 20 spot errors per gel. This amount of error, which might accumulate with each gel to gel comparison used in database construction, could produce an unacceptable degree of uncertainty in gel databases. To address these issues, partial automation of large 2-D gel separations has been undertaken (Nokihara, Morita and Kuriki, 1992; Harrington *et al.*, 1993). Although results are preliminary, spot to spot positional reproducibility in one study was found to be threefold improved over manual methods (Harrington *et al.*, 1993). It should be noted that small 2-D gel formats (50 × 43 mm) have been almost completely automated (Brewer *et al.*, 1986), although these are not generally used for database studies.

#### MICROPREPARATIVE 2-D GEL ELECTROPHORESIS

With the advent of affordable protein microcharacterisation techniques, including N-terminal microsequencing, amino acid analysis, peptide mass fingerprinting, phosphate analysis and monosaccharide compositional analysis, a new challenge for 2-D electrophoresis has been to maintain high resolution and reproducibility but to provide protein in sufficient quantities for chemical analysis (high nanogram to low microgram quantities of protein per spot). This becomes difficult to achieve with very complex samples such as whole bacterial cells, as the initial protein load is divided among 2000 to 4000 protein species. Two approaches are used for producing amounts of material that can be chemically characterised. The first method is to run multiple gels, collect and pool the spots of interest, and subject them to concentration (Ji *et al.*, 1994; Walsh *et al.*, 1995; Rasmussen *et al.*, 1992). In this approach, the concentration process must also act as a purification step to remove accumulated electrophoretic contaminants such as glycine. A more elegant approach has been to exploit the high loading capacity of IPG isoelectric focusing. The high loading capacity of immobilised pH gradients was described early (Ek, Bjellqvist and Righetti, 1983), but has only recently been applied to 2-D electrophoresis (Hanash *et al.*, 1991; Bjellqvist *et al.*, 1993b). Up to 15 mg of protein can be applied to a single gel, yielding microgram quantities of hundreds of protein species. A further benefit of this approach is that proteins present in low abundance, which may not be visualised by lower protein loads, are more likely to be detected. The use of electrophoretic or chromatographic prefractionation techniques (Hochstrasser *et al.*, 1991a; Harrington *et al.*, 1992), followed by high loading of narrow-range IPG separations (Bjellqvist *et al.*, 1993b) provides a likely solution to studies on proteins present in low abundance.

#### Methods of protein detection

There are many means for detecting proteins from 2-D gels. The method used will be dictated by factors including protein load on gel (analytical or preparative), the purpose of the gel (for protein quantitation or for blotting and chemical characterisation), and the sensitivity required. The most common means of protein detection and their applications are shown in Table 1. Most detection methods have drawbacks, for

Table 1: Common stains for 2-D gels or blots and their applications.

Detection Method	Main applications	Unsuitable applications	Sensitivity	References
$^{35}\text{S}$ Met or $^{14}\text{C}$ radiolabelling and fluorimetry or phosphorimaging	Cell lines, cultured organisms	Samples that cannot be labelled	20 ppm of radiolabel in a spot	Ganley and Franza, 1990 Latham, Carre's and Solter, 1993
$^{111}\text{In}$ tinctoria silver	Extremely high sensitivity gel staining	Preparative 2-D, PVDF or NC membranes	10 ng protein on spot or band of gel	Wallace and Saluz, 1992a,b
Silver	Very high sensitivity gel staining, can be more or polychromatic	Preparative 2-D, PVDF or NC membranes	4 ng protein on spot or band of gel	Rahilrud, 1992 Hochstrasser and Merril, 1988
Coomassie blue R-250	Staining of gels, staining of PVDF membranes before protein sequencing	Staining prior to direct mass determination from PVDF, amino acid analysis on PVDF, detection of some glycoproteins	40 ng protein on band or spot of gel	Sirupat <i>et al.</i> , 1994, Gharahdaghi <i>et al.</i> , 1992, Goldberg <i>et al.</i> , 1988, Sanchez <i>et al.</i> , 1992
Colloidal gold	Staining NC membranes, staining PVDF before direct MALDI-TOF	Gels	60% higher than coomassie	Yamaguchi and Asakawa, 1988, Eckerskorn <i>et al.</i> , 1992, Sirupat <i>et al.</i> , 1994
Zinc imidazole	Reverse staining of gels or membranes; may be beneficial in MALDI-TOF of peptides	Where positive image is required	Higher than coomassie	Ortiz <i>et al.</i> , 1992, James <i>et al.</i> , 1993
Ponceau S and amido black	Staining higher protein loads on PVDF, for protein sequencing or amino acid analysis	Staining prior to direct mass determination from PVDF	100 ng protein on band or spot of gel	Sanchez <i>et al.</i> , 1992, Sirupat <i>et al.</i> , 1994, Wilkins <i>et al.</i> , 1995
India ink	Staining of membrane-bound proteins, staining PVDF before direct MALDI-TOF	Gel staining, not quantitative from protein to protein	1-10 ng	Li <i>et al.</i> , 1989, Hughes, Mack and Hamparian, 1988, Sirupat <i>et al.</i> , 1994
Stain-all	Staining to detect glycoproteins or Ca <sup>2+</sup> -binding proteins	General gel staining	100 ng protein on band or spot of gel	Campbell, MacLennan and Jorgensen, 1983, Goldberg <i>et al.</i> , 1988

PVDF = polyvinylidene difluoride; NC = nitrocellulose; MALDI-TOF = matrix assisted laser desorption/ionisation time of flight mass spectrometry.

example, some glycoproteins are not stained by coomassie blue (Goldberg *et al.*, 1988), and many organic dyes are unsuitable for protein detection on PVDF if samples are to be used for direct matrix-assisted laser desorption/ionisation mass spectrometry (Sirupat *et al.*, 1994).

Although most means of protein detection give some indication of the quantities of protein present, in general they cannot be used for global quantitation. This is because

no protein, stain is able consistently to detect proteins over a wide range of concentrations, isoelectric points and amino acid compositions, and with a variety of post-translational modifications (Goldberg *et al.*, 1988; Li *et al.*, 1989). Furthermore, there are large differences in staining pattern when identical gels or blots are subjected to different stains, including amido black, imidazole zinc, india ink, ponceau S, colloidal gold, or coomassie blue (Tovey, Ford and Baldo, 1987; Ortiz *et al.*, 1992). The most common means of quantitating large numbers of proteins in a 2-D gel involves the radiolabelling of protein samples prior to electrophoresis, and protein quantitation based on fluorography and image analysis or liquid scintillation counting (Gurelli, 1989; Celis and Olsen, 1994). However, proteins which do not contain methionine cannot be detected if only [<sup>35</sup>S] methionine is used for labelling. Amino acid analysis of protein spots visualised by other techniques prevents a likely means of protein quantitation for the future.

#### BLOTTING OF PROTEINS TO MEMBRANES

Electrophoretic blotting of proteins from two-dimensional polyacrylamide gels to membranes presents many options for protein identification and microcharacterisation which are not possible when proteins remain in gels. For example, when proteins are blotted to polyvinylidene difluoride (PVDF) membranes, they can be identified by N-terminal sequencing, amino acid analysis, or immunoblotting, or they may be subjected to endoprotease digestion, monosaccharide analysis, phosphate analysis, or direct matrix-assisted laser desorption/ionisation mass spectrometry (Matsudaira, 1987; Wilkins *et al.*, 1995; Jungblut *et al.*, 1994; Sutton *et al.*, 1995; Rasmussen *et al.*, 1994; Weizthandler *et al.*, 1993; Murthy and Iqbal, 1991; Eckerskorn *et al.*, 1992). It is possible to combine some of these procedures on a single protein spot on a PVDF membrane (Packer *et al.*, 1995; Wilkins *et al.*, submitted; Weizthandler *et al.*, 1993). This is useful when minimal amounts of protein are available for analysis. These techniques will be explored in detail later in this review. Notwithstanding the above, there are some disadvantages associated with blotting of proteins to membranes. There is always loss of sample during blotting procedures (Eckerskorn and Lottspeich, 1993), and common protein detection methods are less sensitive or not applicable to membranes (Table 1), presenting difficulties for the analysis of low abundance proteins. Detailed discussion of the merits of available membranes and common blotting techniques can be found elsewhere (Eckerskorn and Lottspeich, 1993; Strupat *et al.*, 1994; Patterson, 1994).

#### 2-D gel analysis, documentation, and proteome databases

Following protein electrophoresis and detection, detailed analysis of gel images is undertaken with computer systems. For proteomic projects, the aim of this analysis is to catalogue all spots from the 2-D gel in a qualitative and if possible quantitative manner, so as to define the number of proteins present and their levels of expression. Reference gel images, constructed from one or more gels, form the basis of two-dimensional gel databases. These databases also contain protein spot identities and

GEL

Att

ph

sca

Cel

res

or r

pul

spo

spe

ass

list

Tab

Gel

ELS

GEL

ME

QUI

TY

Ti

ity

IG

20

im

ma

to

usi

Ch

alt

19

details of their post-translational modifications. 2-D gel databases are beginning to be linked to or integrated with comprehensive protein and nucleic acid databases (Neidhardt *et al.*, 1989; Simpson *et al.*, 1992; Appel *et al.*, 1994), and 'organism' databases, containing DNA sequence data, chromosomal map locations, reference 2-D gels and protein functional information for an organism, are becoming established as genome and proteome projects progress (VanBogelen *et al.*, 1992; Yeast Protein Database cited in Garrels *et al.*, 1994).

#### GEL IMAGE ANALYSIS AND REFERENCE GELS

After 2-D electrophoresis and protein visualisation by staining, fluorography or phosphorimaging, images of gels are digitised for computer analysis by an image scanner, laser densitometer, or charge-coupled device (CCD) camera (Garrels, 1989; Celis *et al.*, 1990a; Urwin and Jackson, 1993). All systems digitise gels with a resolution of 100–200 mm, and can detect a wide range of densities or shading (256 or more 'grey scales'). Following this, gel images are subjected to a series of manipulations to remove vertical and horizontal streaking and background haze, to detect spot positions and boundaries, and to calculate spot intensity (Figure 3). A standard spot (SSP) number, containing vertical and horizontal positional information, is assigned to each detected spot and becomes the protein's reference number. Table 2 lists some notable software packages which process 2-D gel images.

Table 2. Some Software Packages for the Analysis of Gel Images

Gel Image Analysis System	References*
ELSIE I & II	Olsen and Miller, 1988; Wirth <i>et al.</i> , 1991; Wirth <i>et al.</i> , 1993
GELLAB I & II	Wu, Lemkin and Upmum, 1993; Lemkin, Wu and Upmum, 1993; Myrick <i>et al.</i> , 1993
MELANIE I & II	Appel <i>et al.</i> , 1991; Hochstrasser <i>et al.</i> , 1991b
QUEST I & II and PDQUEST	Garrels, 1989; Monardo <i>et al.</i> , 1992; Hohn <i>et al.</i> , 1992; Celis <i>et al.</i> , 1994a,b
TYCHO & KEPLAR	Anderson <i>et al.</i> , 1992; Richardson, Hiron and Anderson, 1992

\* These references are not exhaustive: they include some references of use as well as authors of the system.

As there are difficulties in the electrophoresis of samples with 100% reproducibility, reference gel images are often constructed from many gels of the same sample (Garrels and Franza, 1989; Neidhardt *et al.*, 1989). Since this involves the matching of 2000 to 4000 proteins from one gel to another, it presents a considerable challenge to image analysis systems. Matching of gels is usually initiated by an operator, who manually designates approximately 50 or so prominent spots as 'landmarks' on gels to be cross-matched. Proteins which match are then established around landmarks, using computer-based vector algorithms to extend the matching over the entire gel. Close to 100% of spots from complex samples can be matched by these methods, although different degrees of operator intervention may be required (Olsen and Miller, 1988; Lemkin and Lester, 1989; Garrels, 1989; Myrick *et al.*, 1993).

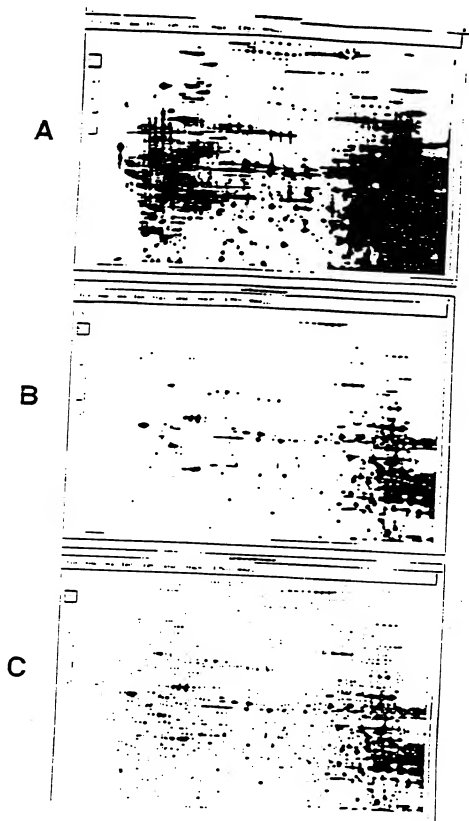


Figure 3. Computer processing of gel images. Shown is a wide pI range 2-D separation of human liver proteins, processed by Melanie software (Appel *et al.*, 1991). (A) Original gel image as captured by laser densitometer. (B) Gel image after processing to remove streaking and background. (C) Outline definition of all spots on the gel.

Ex-  
DQ  
during  
are re-  
obtain  
curves  
MW  
Bogel  
all, 19  
to PNI  
*et al.*  
protein  
amino  
carry  
positio

#### SPOT C

A maj-  
separa  
to dete  
positio  
of maj  
perform  
minut  
scintill  
by rel  
protein  
*et al.*  
Garrel  
be no  
Limit  
not ac  
only e  
an alt  
Myric  
1992  
Wh  
and m  
protei  
their  
transf  
regula  
synth  
(Lath.



## CALCULATION OF PROTEIN ISOELECTRIC POINT AND MOLECULAR WEIGHT

Estimation of the isoelectric point (pI) and molecular weight (MW) of proteins from 2-D gel provides fundamental parameters for each protein, which are also of use during identification procedures (see following section). The pI and MW of proteins are recorded in 2-D gel databases. Accurate estimations of protein pI and MW can be obtained by using 20 or more known proteins on a reference map to construct standard curves of pI and molecular weight, which are then used to calculate estimated pI and MW of unknown proteins (Neidhardt *et al.*, 1989; Garrels and Franza, 1989; Van-Begelen, Hulton and Neidhardt, 1990; Anderson and Anderson, 1991; Anderson *et al.*, 1991; Latham *et al.*, 1992). Alternatively, the MW of individual proteins blotted to PVDF can be determined very accurately by direct mass spectrometry (Eckerskorn *et al.*, 1992). Where immobilised pH gradients are used, the focusing position of proteins allows their pI to be measured within 0.15 units of that calculated from the amino acid sequence (Byellqvist *et al.*, 1993c). It must be noted, however, that proteins carrying post-translational modifications may migrate to unexpected pI or MW positions during electrophoresis (Packer *et al.*, 1993).

## SPOT QUANTITATION AND EXPRESSION ANALYSIS

A major challenge faced in proteome projects is the quantitative analysis of proteins separated by 2-D electrophoresis. The most accurate means of protein quantitation is to determine chemically the amount of each protein present by amino acid composition analysis. However, the current method of choice for quantitative analysis of many proteins is to radiolabel samples with [<sup>35</sup>S] methionine or [<sup>14</sup>C] amino acids, perform the 2-D electrophoresis, and measure protein levels in disintegrations per minute (dpm) or units of optical density. Quantitation is achieved either by liquid scintillation counting, or by gel image analysis where spot densities are quantitated by reference to gel calibration strips containing known amounts of radiolabelled protein or against the integrated optical density of all spots visualised (Vandekerckhove *et al.*, 1990; Celis *et al.*, 1990b; Celis and Olsen, 1994; Garrels, 1989; Latham, Garrels and Solier, 1993; Fey *et al.*, 1994). All approaches effectively allow spots to be normalised against the total disintegrations per minute loaded onto the gel. Limitations that remain with radiolabelling methods are that absolute quantitation is not achieved because all proteins have varying amounts of any amino acid, and that only easily labelled samples can be investigated. Quantitative silver staining presents an alternative (Girometti *et al.*, 1991; Harrington *et al.*, 1992; Rodriguez *et al.*, 1993; Myrick *et al.*, 1993), which when undertaken with [<sup>35</sup>S]thiourea (Wallace and Saluz, 1992) allows of extremely high sensitivity.

When protein spots from samples prepared under different conditions are quantitated and matched from gel to gel, it becomes possible to examine changes and patterns in protein expression. Large scale investigation of up- and down-regulation of proteins, their appearance and disappearance, can be undertaken. For example, simian virus 40 transformed human keratinocytes were shown to have 177 up-regulated and 58 down-regulated proteins compared to normal keratinocytes (Celis and Olsen, 1994); detailed synthesis profiles of 1200 proteins have been established in 1 to 4 cell mouse embryos (Latham *et al.*, 1991, 1992); and 4 proteins out of 1971 were found to be markers for

cadmium toxicity in urinary proteins (Myrick *et al.*, 1993). Complex global changes in protein expression as a result of gene disruptions have also been investigated (S. Fey and P. Moss-Larsen, Personal communication). Impressively, large gel sets showing protein expression under different conditions can be globally investigated using statistical methods that find groups of related objects within a set. For example, the REF52 rat cell line database, consisting of 79 gels from 12 experimental groups where each gel contains quantitative data for 1600 cross-matched proteins, has been analysed by cluster analysis (Garrels *et al.*, 1990). This revealed clusters of proteins that, for example, were induced or repressed similarly under simian virus 40 or adenovirus transformation, suggesting a common mechanism. Protein groups that were induced or repressed during culture growth to confluence were also found. It is obvious that the potential for investigation of cellular control mechanisms by these approaches is immense. It is equally clear that investigations of gene expression of this scale are currently technically impossible using nucleic-acid based techniques.

Table 3. Some proteomic databases and their special features

Proteomic database	Special features	References
<i>E. coli</i> gene-protein database	Gel spots linked with GenBank and Kohara clones; quantitative spot measurements under different growth conditions	VanBogelen and Neidhardt, 1991; VanBogelen <i>et al.</i> , 1992
Human heart database	Identification of disease markers; two separate databases have been established	Baker <i>et al.</i> , 1992; Corbett <i>et al.</i> , 1992b
Human keratinocyte database	Extensive identifications; quantitative spot measurements of transformed cells; identification of disease markers	Jungblut <i>et al.</i> , 1992; Celis <i>et al.</i> , 1993b; Celis <i>et al.</i> , 1993; Celis and Olsen, 1994
Mouse embryo database	Quantitative spot measurements through 1 to 4 cell stage	Latham <i>et al.</i> , 1991; Latham <i>et al.</i> , 1992
Mouse liver database (Arginine Protein Mapping Group)	Documented changes due to exposure to ionizing radiation and toxic chemicals	Grimm, Taylor and Tøllaksen, 1992
Rat liver epithelial database	Deviated subcellular fractionation studies	Wirth <i>et al.</i> , 1991; Wirth <i>et al.</i> , 1992
Rat liver database	Extensive studies on regulation of proteins by drugs and toxic agents	Anderson and Anderson, 1991; Anderson <i>et al.</i> , 1992; Richardson, Horn and Anderson, 1992
REF 52 rat cell line database	Accessible via World Wide Web; quantitative spot measurements under different conditions	Garrels and Franza 1989; Bruni <i>et al.</i> , 1992
SWISS-2DPAGE containing human reference maps	Accessible via World Wide Web; completely integrated with SWISS-PROT and SWISS-3DIMAGE	Appel <i>et al.</i> , 1993; Hochstrasser <i>et al.</i> , 1992; Hughes <i>et al.</i> , 1993; Golzar <i>et al.</i> , 1993
Yeast Protein Database (YPD) and Yeast Electrophoretic Protein Database (YEPD)	Completely crossreferenced organism database; YPD has extensive information on over 3500 proteins; YEPD has many identifications	Garrels <i>et al.</i> , 1994

Protein  
protein  
informa  
2-D ge  
subset  
of refe  
should  
Macini  
the are  
annota  
sequen  
One  
SWISS  
1994;  
feature  
2DPA

Table 4  
All three  
expans  
Inform

Annex

Cross-  
Refer  
Data

Other

## FEATURES OF PROTEOME DATABASES

Proteome projects rely heavily on computer databases to store information about all proteins expressed by an organism. 'Proteome databases' should contain detailed information of proteins already characterised elsewhere, as well as protein data from 2-D gels such as apparent pI and MW, expression level under different conditions, subcellular localisation, and information on post-translational modifications. Images of reference 2-D gels, showing protein SSP numbers and protein identifications, should also be included. Ideally, proteome databases should be accessible with Macintosh or IBM personal computers and easy to use. Some proteome databases and the areas they cover are listed in Table 3. Databases range from collections of unannotated gels to large databases of images integrated with protein and nucleic acid sequence banks.

One example of an integrated proteome database is the suite of SWISS-PROT, SWISS-2DPAGE and SWISS-3DIMAGE databases (Appel *et al.*, 1993; Appel *et al.*, 1994; Appel, Bairoch and Hochstrasser, 1994; Bairoch and Boeckmann, 1994). The features of these three databases are listed in Table 4. SWISS-PROT, SWISS-2DPAGE and SWISS-3DIMAGE are accessible through the World Wide Web

Table 4: The SWISS-PROT, SWISS-2DPAGE and SWISS-3DIMAGE suite of crosslinked databases. All three databases are accessible through the World Wide Web, at URL address: <http://expasy.baug.ch/>

	SWISS-PROT	SWISS-2DPAGE	SWISS-3DIMAGE
Information	Text entries of sequence data. Chain information. Taxonomic data, 36, 303 entries in Release 29	2-D gel images of: human liver, plasma, HepG2, HepG2 secreted proteins, red blood cell, lymphoma, cerebrospinal fluid, macrophage like cell line, erythroleukemia cell, platelet	Collection of 341 2-D images of proteins
Annotations	Protein function. Post translational modifications. Domains. Secondary structure. Quaternary structure. Diseases associated with protein. Sequence conflicts	Gel images where protein is found. How protein identified. Protein pI and MW, protein number, normal and pathological variants	All annotation is available in SWISS- PROT
Cross- Referenced Databases	SWISS-2DPAGE SWISS-3DIMAGE EMBL, PIR, PDB, OMIM, PROSITE, Medline, Flybase, GCRDB, MaizeDB, WormPep, DictyDB	SWISS-PROT and all other databases accessible through SWISS-PROT	SWISS-PROT and all other databases accessible through SWISS-PROT
Other Features	Navigation to other SWISS databases achieved by selecting entries with computer mouse	Gel images show position of identified proteins, or region of gel where protein should appear	Mono and stereo images available. Images can be transferred to local computer image viewing programs

(Berners-Lee *et al.*, 1992), allowing any computer connected to the internet to access the stored information and images. Navigation within and between the three databases is seamless, as all potential crosslinks are highlighted as hypertext on the display and can be selected with a computer mouse. From these databases, detailed information about a protein, including amino acid sequence and known post-translational modifications, can be obtained, the precise protein spot it corresponds to on a reference gel image can be viewed if known, and the 3-D structure of the molecule can be seen if available. References to nucleic acid and other databases are also given to provide access to information stored elsewhere.

Organism databases, containing detailed protein and nucleic acid information about a species, are becoming common as genome and proteome projects progress. These differ from nucleic acid or protein sequence databases like GenBank or SWISS-PROT because they are image based, and contain information about chromosomal map positions, transcription of genes, and protein expression patterns. The *Escherichia coli* gene-protein database (VanBogelen, Hutton and Neidhardt, 1990; VanBogelen and Neidhardt, 1991; VanBogelen *et al.*, 1992), known as the ECO2DBASE, is one example. It contains gene and protein names, 2-D gel spot information (including pI and MW estimates, and spot identification), genetic information (GenBank or EMBL codes, chromosomal location, location on Kohara clones (Kohara, Akiyama, and Isono, 1987), transcription direction of genes), and protein regulatory information (level of protein expression under different growth regimes, member of regulon or stimulon). All entries in the ECO2DBASE are also cross-referenced to the SWISS-PROT database (Bairoch and Boeckmann, 1994). It is anticipated that organism databases will soon become a standard means of storing all available information about a particular species. However there is currently no consistent manner in which organism databases are assembled, which may hamper comparisons in the future.

### Identification and characterisation of proteins from 2-D gels

The number of proteins identified on a 2-D reference map determines its usefulness as a research and reference tool. As most reference maps have only a small proportion of proteins identified, a major aim of current proteome projects is to screen many proteins from 2-D maps, in order to define them as 'known' in current nucleic acid and protein databases, or as 'unknown'. Protein identification assists in confirmation of DNA open reading frames, and provides focus for DNA sequencing projects, and protein characterisation efforts by pointing to proteins that are novel. Since there may be 3000–10000 proteins from a single 2-D map that require identification, the challenge in protein screening is to identify proteins quickly, with a minimum of cost and effort.

Traditionally, proteins from 2-D gels have been identified by techniques such as immunoblotting, N-terminal microsequencing, internal peptide sequencing, comigration of unknown proteins with known proteins, or by overexpression of homologous genes of interest in the organism under study (Matsudaira, 1987; Rosenfeld *et al.*, 1992; VanBogelen *et al.*, 1992; Celis *et al.*, 1993; Honore *et al.*, 1993; Garrels *et al.*, 1994). Whilst these techniques are powerful identification tools, they are too expensive or time and labour intensive to use in mass screening programs. A hierarchical approach to mass protein identification has been recently suggested as an

alien  
use c  
mass  
slow  
the c  
of th  
nucl  
cons  
tech  
ident

PROT

Ther  
ident  
This  
to id  
The  
radi  
al.,  
chro  
1981  
1991  
phor  
radi

**Table 5:** Hierarchical analysis for mass screening of 2-D separated proteins: blotted to membrane. Rapid and inexpensive techniques are used as a first step in protein identification, and then more expensive techniques are then used if necessary. Table modified from Wasinger *et al.*, 1995.

Order	Identification technique	Reference
1	Amino acid analysis	Jungblut <i>et al.</i> , 1992; Sliwa, 1997; Hohmann, Houthuys and Sander, 1994; Jungblut <i>et al.</i> , 1992; Wilkins <i>et al.</i> , 1995
2	Amino acid analysis with N-terminal sequence tag	Wilkins <i>et al.</i> , submitted
3	Peptide-mass fingerprinting	Henzel <i>et al.</i> , 1993; Pappin, Hourup and Bidaux, 1993; Jones <i>et al.</i> , 1993; Mann, Hourup and Roepstorff, 1993; Yates <i>et al.</i> , 1993; Altmann <i>et al.</i> , 1992; Sutton <i>et al.</i> , 1995
4	Combination of amino acid analysis and peptide mass fingerprinting	Cordwell <i>et al.</i> , 1995; Wasinger <i>et al.</i> , 1995
5	Mass spectrometry sequence tag	Mann and Wilm, 1994
6	Extensive N-terminal Edman microsequencing	Marudaira, 1987
7	Internal peptide Edman microsequencing	Rosenfeld <i>et al.</i> , 1992; Hellman <i>et al.</i> , 1995
8	Microsequencing by mass spectrometry, electrospray ionisation, post-source decay MALDI-TOF	Johnson and Walsh, 1992
9	Ladder sequencing	Barile-Jones <i>et al.</i> , 1992

alternative to traditional approaches (Table 5; Wasinger *et al.*, 1995). This involves the use of rapid and cheap identification tools such as amino acid analysis and peptide mass fingerprinting as first steps in protein identification, followed by the use of slower, more expensive and time consuming identification procedures if necessary. In the construction of this hierarchy the analysis time, cost per sample and the complexity of the data created has been considered, as whilst some techniques require little machine time per sample, the analysis of data can be quite involved and time consuming. Amino acid analysis and peptide mass-fingerprinting based identification techniques in the hierarchy are discussed in detail below. For review of other protein identification techniques in Table 5, see Patterson (1994) and Mann (1995).

#### PROTEIN IDENTIFICATION BY AMINO ACID COMPOSITION

There has been a revival of interest in the use of amino acid composition for identification of proteins from 2-D gels after early work by Eckerskorn *et al.* (1988). This technique uses a protein's idiosyncratic amino acid composition profile in order to identify it by comparison with theoretical compositions of proteins in databases. The amino acid composition of proteins can be determined by differential metabolic radiolabelling and quantitative autoradiography after 2-D electrophoresis (Garrels *et al.*, 1994; Frey *et al.*, 1994), or by acid hydrolysis of membrane-blotted proteins and chromatographic analysis of the resulting amino acid mixture (Eckerskorn *et al.*, 1988; Tous *et al.*, 1989; Gharahdaghi *et al.*, 1992; Jungblut *et al.*, 1992; Wilkins *et al.*, 1995). As differential metabolic labelling experiments require X-ray film or phosphor-image plate exposures of up to 140 days, and can only be undertaken with easily radiolabelled samples, the technique is not as rapid or widely applicable as chromato-

## Spot: ECOLI-BIM

\*\*\*\*\*

## Composition:

Asx: 10.2 Glx: 10.4 Ser: 5.7 His: 0.7  
 Gly: 5.4 Thr: 1.8 Ala: 6.7 Pro: 7.9  
 Tyr: 0.3 Asp: 5.0 Val: 8.0 Met: 0.3  
 Ile: 5.9 Leu: 8.0 Phe: 13.3 Lys: 4.4

pI estimate: 6.89 Range searched: ( 6.64, 7.14)  
 Mw estimate: 16800 Range searched: (13440, 20160)

Closest: SWISS-PROT entries for the species ECOLI matched by AA composition:

Rank	Score	Protein	pI	Mw	Description
1	24	PYR_ECOLI	6.84	16989	ASPARTATE CARBOXYLTRANSFERASE
2	39	CSMA_ECOLI	6.32	36359	PANTOTHENATE KINASE (EC 2.7.1.33)
3	40	MTA_ECOLI	5.06	35723	MONOISOMERIC O-SUCCINYLTRANSFERASE
4	42	DAD_ECOLI	5.52	57812	TRANSCRIPTIONAL ACTIVATOR CADC.
5	43	HLYS_ECOLI	8.36	19769	HEMOLYSIN C. PLASMED.

Closest: SWISS-PROT entries for ECOLI with pI and Mw values in specified range:

Rank	Score	Protein	pI	Mw	Description
1	24	PYR_ECOLI	6.84	16989	ASPARTATE CARBOXYLTRANSFERASE
2	122	TRAF_ECOLI	6.73	17921	TRAF PROTEIN.
3	122	YAPB_ECOLI	6.79	19028	HYPOTHETICAL LIPOPROTEIN YAPB.
4	140	YFSP_ECOLI	6.83	14945	HYPOTHETICAL 14.9 KD PROTEIN IN GRPE
5	142	YAPB_ECOLI	7.06	14726	HYPOTHETICAL PROTEIN IN BETT 3' REGION

Figure 4. Computer printout from EAPASy server where the empirical amino acid composition, estimated pI and MW of a protein from a 2-D reference map of *E. coli* were matched against all entries in SWISS-PROT for *E. coli*. The correct identification, aspartate carboxyltransferase, is shown in bold. Low scores indicate a good match. Note how matching within a defined pI and MW range (lower set of proteins) has greatly increased the score difference between the first and second ranking proteins. This score difference gives high confidence in the identification, and is only observed where the top ranking protein is the correct identification (Wilkins *et al.*, 1995).

graphically-based analysis. Proteins blotted in PVDF membranes can be hydrolysed in 1 h at 155°C, amino acids extracted in a single brief step, and each sample automatically derivatised and separated by chromatography in under 40 minutes (Wilkins *et al.*, 1995; Ou *et al.*, 1995). In this manner, one operator can routinely analyse 100 proteins per week on one HPLC unit. This technology lends itself to automation, and it is anticipated that instruments with even greater sample throughput will be developed. When proteins have been prepared by micropreparative 2-D electrophoresis (Hanash *et al.*, 1991; Bjellqvist *et al.*, 1993b), blotted to a PVDF membrane and stained with amido black, any visible protein spot is of sufficient quantity for amino acid analysis (Cordwell *et al.*, 1995; Waxinger *et al.*, 1995; Wilkins *et al.*, 1995).

After the amino acid composition of a protein has been determined, computer programs are used to match it against the calculated compositions of proteins in databases (Eckerskorn *et al.*, 1988; Sibbald, Sommerfeldt and Argos, 1991; Jungblut *et al.*, 1992; Shaw, 1993; Hobohm, Houthaeve and Sander, 1994; Wilkins *et al.*, 1995). Matching is usually done with only 15 or 16 amino acids, as cysteine and

Figure 4  
 same as  
 acid com  
 PROT: f  
 for this  
 large se  
 the over  
 protein

trypto  
 to thei  
 The co  
 a score  
 restric  
 1994:  
*et al.*  
 match  
 in Fig  
 refere  
 runny  
 lymph  
*et al.*

PROT:  
 SEQU  
 When

## Spot ECOLI-ASP

\*\*\*\*\*

## Composition:

Asx: 9.4 Glx: 10.8 Ser: 4.1 His: 2.7  
 Gly: 10.2 Thr: 2.8 Ala: 11.9 Pro: 2.2  
 Tyr: 6.0 Arg: 3.7 Val: 9.5 Met: 0.6  
 Ile: 5.1 Leu: 8.2 Phe: 3.2 Lys: 6.9

pI estimate: 5.99 Range searched: ( 5.74, 6.24)  
 Mw estimate: 45000 Range searched: (36000, 54000)

Closest SWISS-PROT entries for ECOLI with pI and Mw values in specified range:

Rank	Score	Protein	pI	Mw	N-terminal Seq.
1	21	GLXA_ECOLI	6.03	45316	M L K R E
2	22	YHGB_ECOLI	5.86	36502	H S M I K
3	38	GART_ECOLI	5.78	45774	H S N S K
4	44	YHNS_ECOLI	5.86	48018	H R I K Y
5	45	DHKA_ECOLI	5.98	48581	H D Q T Y
6	46	ARGD_ECOLI	5.79	43765	H A I E Q
7	46	MYRA_ECOLI	5.78	37851	H N H S L
8	47	GLMY_ECOLI	5.98	49162	H L N R A
9	47	AKRA_ECOLI	5.85	43290	H S S K L
10	50	YHSH_ECOLI	6.01	37064	H E S R I

Figure 5. A PVDF protein spot from an *E. coli* 2-D reference map was sequenced for 4 cycles, and the same sample then subjected to amino acid analysis. The N-terminal sequence was M L K R. When the amino acid composition of the spot, as well as estimated pI and MW, were matched against all entries in SWISS-PROT for *E. coli*, the above list of best matches was produced. N-terminal sequences are from SWISS-PROT for those entries. The top ranking identification of serine hydroxymethyltransferase (bold) did not show a large score difference between the first and second ranking proteins, giving little confidence in this being the correct protein identification. However, the sequence tag M L K R confirmed the identity of the protein as serine hydroxymethyltransferase.

tryptophan are destroyed during hydrolysis, asparagine and glutamine are deamidated to their corresponding acids, and proline is not quantitated in some analysis systems. The computer programs produce a list of best matching proteins, which are ranked by a score that indicates the match quality. Some programs allow matching to be restricted to specific 'windows' of MW and pI (Hobohm, Houtchaev and Sander, 1994; Wilkins *et al.*, 1995), and to protein database entries for one species (Jungblut *et al.*, 1992; Wilkins *et al.*, 1995). The use of such restrictions increases the power of matching. An example of protein identification by amino acid composition is shown in Figure 4. To date, amino acid composition has been used to identify proteins from reference maps of *Spiroplasma melliferum*, *Mycoplasma genitalium*, *E. coli*, *Saccharomyces cerevisiae*, *Dictyostelium discoideum*, human sera, human heart, human lymphocyte, and mouse brain (Cordwell *et al.*, 1995; Wasinger *et al.*, 1995; Wilkins *et al.*, 1995; Jungblut *et al.*, 1992, 1994; Garrels *et al.*, 1994; Frey *et al.*, 1994).

# PROTEIN IDENTIFICATION BY AMINO ACID COMPOSITION AND N-TERMINAL SEQUENCE TAG

When samples from 2-D gels are not unambiguously identified by amino acid

composition, pI and MW, often the correct identification of that protein is amongst the top rankings of the list (Hobohm, Houthaeve and Sander, 1994; Cordwell *et al.*, 1995; Wilkins *et al.*, 1995). Taking advantage of this observation, we have used the mass spectrometry 'sequence tag' concept (Mann and Willen, 1992) in developing a combined Edman degradation and amino acid analysis approach to protein identification (Wilkins *et al.*, submitted). This involves the N-terminal sequencing of PVDF-blotted proteins by Edman degradation for 3 or 4 cycles to create a 'sequence tag', following which the same sample is used for amino acid analysis. As only a few amino acids are removed from the protein, its composition is not significantly altered. Furthermore, since only a small amount of protein sequence is required, fast but low repetitive yield Edman degradation cycles can be used. Modifications to current procedures should allow 3 cycles to be completed in 1 h, thereby allowing the screening of 100 or more proteins per week on one automated, multi-cartridge sequencer. Amino acid composition, pI and MW of proteins are matched against databases as described above, and N-terminal sequences of best matching proteins are checked with the 'sequence tag' to confirm the protein identity (Figure 5). This technique will be less useful when proteins are N-terminally blocked, but as only a few N-terminal amino acids are susceptible to the acetyl, formyl, or pyroglutamyl modifications that cause blockage, this may itself provide useful information for sequence tag identification. A strength of N-terminal sequence tag and amino acid composition protein identification is that data generated are quickly and easily interpreted.

#### PROTEIN IDENTIFICATION BY PEPTIDE MASS FINGERPRINTING

Techniques for the identification of proteins by peptide mass fingerprinting have recently been described (Henzel *et al.*, 1993; Pappin, Hojrup and Bleasby, 1993; James *et al.*, 1993; Mann, Hojrup and Roepstorff, 1993; Yates *et al.*, 1993; Moritz *et al.*, 1994; Sutton *et al.*, 1995). This involves the generation of peptides from proteins using residue-specific enzymes, the determination of peptide masses, and the matching of these masses against theoretical peptide libraries generated from protein sequence databases. As proteins have different amino acid sequences, their peptides should produce characteristic 'fingerprints'.

The first step of peptide mass fingerprinting is protein digestion. Proteins within the gel matrix or bound to PVDF can be enzymatically digested *in situ*, although *in situ* gel digests are reported to produce more enzyme autodigestion products, which complicate subsequent peptide mass analysis (James *et al.*, 1993; Rasmussen *et al.*, 1994; Moritz *et al.*, 1994). The enzyme of choice for digestion is currently trypsin (of modified sequencing grade), but other enzymes (Lys-C or *S. aureus* V8 protease) have also been used (Pappin, Hojrup and Bleasby, 1993). To maximise the number of peptides obtained, it is desirable for protein samples to be reduced and alkylated prior to digestion (Moritz *et al.*, 1994; Henzel *et al.*, 1993). This ensures that all disulfide bonds of the protein are broken, and produces protein conformations that are more amenable to digestion. Surprisingly, chemical digestion methods such as cyanogen bromide (methionine specific), formic acid (aspartic acid specific), and 2-(2'-nitrophenyl)sulfonyl-3-methyl-3-bromoindolenine (tryptophan specific) have not been explored as means of peptide production for mass fingerprinting, even though they are rapid and may circumvent some problems associated with enzyme digestions.



(Nikodem and Freyco, 1979; Crimmins *et al.*, 1990; Vanfleteren *et al.*, 1992).

After proteins are digested, peptide masses are determined by mass spectrometry. Direct analysis of peptide mixtures can be achieved by electrospray ionisation mass spectrometry, plasma desorption mass spectrometry, or matrix assisted laser desorption ionization (MALDI) mass spectrometry techniques. MALDI is preferable because of its higher sensitivity and greater tolerance to contaminating substances from 2-D gels (Janies *et al.*, 1993; Mortz *et al.*, 1994; Pappin, Hojrup and Bleasby, 1993). Furthermore, recent modifications to sample preparation methods have largely solved early difficulties experienced with the calibration of MALDI spectra (Mortz *et al.*, 1994; Vorm and Mann, 1992; Vorm, Roepstorff and Mann, 1994). The high sensitivity of mass spectrometry allows a small fraction of a digest of a 1 µg protein spot to be used for analysis, and analysis itself is complete in a few minutes.

A major challenge associated with peptide mass fingerprinting is data interpretation prior to computer matching against libraries of theoretical peptide digests. Spectra must be examined carefully to determine which peaks represent peptide masses of interest, as there are often enzyme autodigestion products and contaminating substances present (Henzel *et al.*, 1993; Mortz *et al.*, 1994; Rasmussen *et al.*, 1994). Furthermore, if protein alkylation and reduction has not been undertaken prior to protein digestion, peptide sequence coverage may be poor (40% to 70%), with some masses present representing disulfide bonded peptides originally present in the protein (Mortz *et al.*, 1994). For eukaryotes, a serious issue is the alteration of peptide masses by the presence of post-translational modifications (Table 6). The mass of the unmodified peptide alone can be very difficult to determine. Two artifactual modifications introduced by electrophoresis, an acrylamide adduct to cysteine and the oxidation of methionine, are also known to alter peptide masses (Le Maire *et al.*, 1993; Hess *et al.*, 1993).

Table 6: Masses of some common post-translational modifications. Peptides carrying post-translational modifications complicate data analysis for peptide mass fingerprinting protein identification. This is especially so for protein glycosylation, which involves many different combinations of the hexosamines, hexoses, deoxyhexoses, and sialic acid.

Post-translational modification	Mass change
Acetylation	
* Acrylamide adduct to cysteine	+ 22.04
Carboxylation at Asp or Glu	+ 71.00
Hydroxylation at Asn or Gln	+ 24.01
Disulfide bond formation	+ 0.98
Deoxyhexoses (Fuc)	+ 2.02
Formylation	+ 146.14
Hexosamines (GlcN, GalN)	+ 28.01
Hexoses (Glc, Gal, Man)	+ 161.16
Hydroxylation	+ 162.14
N-acetylhexosamines (GlcNAc, GalNAc)	+ 16.00
* Oxidation of Met	+ 203.19
Phosphorylation	+ 16.00
Pre-glutamic acid formed from Gln	+ 79.98
Sialic acid (NeuNAc)	+ 17.03
Sulfation	+ 291.26
	+ 80.06

Table modified from Finnigan LASERMAT application data sheet 4.

\* Asterisk \* shows modifications that can arise artifactually from the 2-D electrophoresis process.

A number of computer programs are available for matching peptide masses against databases (reviewed in Cottrell, 1994). Matching is usually undertaken in an interactive manner, whereby peaks of mass 500–3000 Da are selected and matched under various search parameters including MW of protein, mass accuracy of peptides, and number of missed enzyme cleavages allowed (Henzel *et al.*, 1993; Mortz *et al.*, 1994; Rasmussen *et al.*, 1994). The correct protein identity is the protein which has the most peptide masses in common with the unknown sample. Identities have been established with as few as three peptides, but unambiguous identification is thought to require a mass spectrometric map covering most peptides of the protein (Mortz *et al.*, 1994; Yates *et al.*, 1993). To date, peptide mass fingerprinting of proteins has been undertaken from the human myocardial protein and keratinocyte maps, from an *E. coli* 2-D gel, and from reference maps of *Spiroplasma mellithron* and *Mycoplasma genitalium* (Sutton *et al.*, 1995; Rasmussen *et al.*, 1994; Henzel *et al.*, 1993; Cordwell *et al.*, 1995; Wasinger *et al.*, 1995), although the technique is most powerful when used in combination with another protein identification technique (Rasmussen *et al.*, 1994; Cordwell *et al.*, 1995).

#### MASS SPECTROMETRY SEQUENCE TAGGING

An extension of peptide mass fingerprinting has recently been described, called peptide sequence tagging (Mann and Wilm, 1994; Mann, 1995). This uses tandem mass spectrometry (MS/MS) to initially determine the mass of peptides, then subject them to fragmentation by collision with a gas, and finally determine the mass of fragments. The resulting spectra gives information about a peptide's amino acid sequence. The fragmentation masses of peptides can rarely be used to assign a complete sequence, but it usually allows a short 'sequence tag' of 2 or 3 amino acids to be determined. This sequence tag and the original peptide mass is matched by computer against a database, providing a likely identity of the peptide and the protein it came from. The major drawback for this technique as a mass screening tool is the complexity of the mass data generated and the high level of expertise required for its interpretation. Nevertheless, it represents a useful new protein identification method which greatly increases the power of peptide mass fingerprinting protein identification.

#### Cross-species protein identification

Protein sequence databases continue to grow at a rapid rate, yet it is not widely appreciated that close to 90% of all information contained in current protein databases comes from only 10 species (A. Bairoch, Pers. Comm.). Fortunately, this information can be used to study proteomes of organisms that are poorly defined at the molecular level, via 2-D electrophoresis and 'cross-species' protein identification (Cordwell *et al.*, 1995; Wasinger *et al.*, 1995). This approach allows proteins from reference maps of many different species to be identified without the need for the corresponding genes to be cloned and sequenced. This is particularly true for 'housekeeping' proteins, such as enzymes involved in glycolysis, DNA manipulation and protein manufacture, which are highly conserved across species boundaries. Proteins that cannot be identified across species boundaries can then become the focus of further protein characterisation and DNA sequencing efforts.

A)

Protein APAL\_HUMAN

\*\*\*\*\*

Asx: 8.4 Cys: 19.3 Ser: 6.3 His: 1.3  
 Gly: 4.2 Thr: 4.3 Ala: 8.0 Pro: 4.2  
 Tyr: 1.9 Trp: 6.7 Val: 5.5 Met: 1.3  
 Ile: 6.0 Leu: 15.5 Phe: 2.5 Lys: 8.8

PI Range: no range specified

MW Range: no range specified

The closest SWISS-PROT entries are:

Rank	Score	Protein	pI	MW	Description
1	5	APAL_HUMAN	5.27	28078	APOLIPOPROTEIN A-I.
2	4	APAL_HAIFA	5.43	28005	APOLIPOPROTEIN A-I.
3	12	APAL_RABET	5.15	27836	APOLIPOPROTEIN A-I.
4	14	APAL_BOTIN	5.36	27549	APOLIPOPROTEIN A-I.
5	14	APAL_CANTA	5.10	27467	APOLIPOPROTEIN A-I.
6	18	APAL_MCUSE	5.42	27922	APOLIPOPROTEIN A-I.
7	26	APAL_PIS	5.19	27598	APOLIPOPROTEIN A-I.
8	27	APAL_CWSTX	5.26	27966	APOLIPOPROTEIN A-I.
9	37	DYNA_HCHST	5.44	117742	DYNACTIN, 117 KD ISOFORM.
10	39	APAL_HUMAN	5.18	43374	APOLIPOPROTEIN A-IV.

B)

Reagent: Trypsin MW filter: 10k

Scan using fragment mws of:

1953 1913 1731 1613 1401 1387  
 1201 1283 1252 1235 1231 1215  
 1101 896 873 831 813 781  
 712 704

No. of database entries scanned = 72018

1	APAL_HUMAN	APOLIPOPROTEIN A-I (APO-AI) - HOMO SAPIENS
2	APAL_HAIFA	APOLIPOPROTEIN A-I (APO-AI) - HOMO SAPIENS
3	APAL_PAPHA	APOLIPOPROTEIN A-I (APO-AI) - PAPIO HAMADRYAS
4	B41845	orf B - Treponema denticola
5	APAL_CANTA	APOLIPOPROTEIN A-I (APO-AI) - CANIS FAMILIARIS (DOG)
6	S10947	hypothetical protein - Azotobacter vinelandii
7	MS21_PEA	CHLOROPLAST HEAT SHOCK PROTEIN PRECURSOR - PISUM SATIVUM
8	S20724	Tropomyosin - African clawed frog
9	HIV1354	HIV1354 premature term. at 793 - Human immunodeficiency
10	TRIT_ECOLI	TRAP PROTEIN - ESCHERICHIA COLI

Figure 6. Theoretical cross-species matching of human apolipoprotein A-I by amino acid composition and tryptic peptides. When an unknown protein is analysed, new ranking proteins from both techniques can be compared. If the same protein type is observed in both lists, there is high confidence in this being the identity of the unknown molecule (Cordwell *et al.*, 1995). (A) Output of ExPASy server (Appel, Bairoch and Hochstrasser, 1994) where the true amino acid composition of apolipoprotein A-I was matched against all entries in the SWISS-PROT database, without pI or MW windows. Seven of the top 10 matching proteins were apolipoprotein A-I of different species. (B) Output of MOWSE peptide mass fingerprinting program (Pappin, Horiup and Bleasby, 1993) where true tryptic peptides of human apolipoprotein A-I were matched against the OWL database, using MW window of 10k. Four of the top ten matching proteins were apolipoprotein A-I from different species.

Rapid cross-species identification of proteins from 2-D reference maps can be undertaken with amino acid composition or peptide mass fingerprinting methods (Figure 6), but these techniques alone may not identify proteins unambiguously when phylogenetic cross-species distances are great or analysis data is of poor quality (Yates *et al.*, 1993; Shaw, 1993; Cordwell *et al.*, 1995). However, very high confidence in protein identities can be achieved when lists of best-matching proteins generated by both techniques are compared (Cordwell *et al.*, 1995; Wasinger *et al.*, 1995). The correct identification is found when the same protein is ranked highly in lists of best matches generated by both techniques. This method has allowed approximately 120 proteins from the reference map of the mollicute *Spiroplasma melliferum*, representing approximately one quarter of the proteome, to be confidently identified by reference to protein information from other species (S. Cordwell, Personal Communication). When cross-species protein identification is to be undertaken, it should be noted that the molecular weight of a protein type across species is usually highly conserved, but that protein pI can vary by more than 2 units (Cordwell *et al.*, 1995). Accurate molecular weight determination by direct mass spectrometry of proteins blotted to PVDF (Eckerskorn *et al.*, 1992) should therefore be a useful additional parameter for cross-species protein identification.

#### CHARACTERISATION OF POST-TRANSLATIONAL MODIFICATIONS

Many proteins are modified after translation. Such post-translational modifications, including glycosylation, phosphorylation, and sulfation (see Table 6), are usually necessary for protein function or stability. Some abnormal modifications are associated with disease (Duthel and Revol, 1993; Ghosh *et al.*, 1993; Yamashita *et al.*, 1993). In proteome studies, post-translational modifications can be examined on all proteins present, or on individual spots. Studies on all proteins provide an indication of which proteins may carry a certain type of modification. For example, 2-D gel analysis of cell cultures grown in the presence of [<sup>3</sup>H] mannose or [<sup>32</sup>P] phosphate gives an indication of which proteins carry glycans containing mannose, and which proteins are phosphorylated (Garrel and Franza, 1989). Lectin binding studies of 2-D gels blotted to PVDF or nitrocellulose provide information on the saccharides, if any, that are carried by proteins present (Gravel *et al.*, 1994).

When individual proteins of interest carrying post-translational modifications have been found, micropreparative 2-D electrophoresis can be used to purify them in microgram quantities (Hanash *et al.*, 1991; Bjellqvist *et al.*, 1993b). If protein isoforms of similar MW and pI are to be studied, focusing with narrow range pI gradients (1 pH unit) can provide greater separation and resolution. After electrophoresis, the type and degree of protein phosphorylation can be investigated (Murthy and Iqbal, 1991; Gold *et al.*, 1994), monosaccharide composition can be determined (Weitzhandler *et al.*, 1993; Packer *et al.*, 1995), and the structure and exact site of glycoamino acids can be investigated by either Edman degradation based techniques or by mass spectrometry (Pisano *et al.*, 1993; Huberty *et al.*, 1993; Carr, Huddleston and Bean, 1993). With further development of rapid techniques, investigation of phosphorylation and monosaccharides by chromatographic or mass spectrometric means is likely to become a routine step in the characterisation of post-translational modifications of proteins from reference maps.

The  
Non-  
trans-  
Acid  
mole  
each  
prote  
genom  
side o  
comp  
thru  
plasm  
Wash  
maps  
specie  
and si

Table  
PROT  
referen  
from B  
1996  
Species

Mycop.  
Escherri  
Sacchari  
Dietrich  
Armbul.  
Caenorh  
Hemo

The  
under  
hecau  
hundr  
estim:  
to tiss  
protei  
electr  
ism c  
accel  
are ur  
post-i  
differ  
useful

## The status of proteome projects

Many technical aspects of proteome research have already been discussed in this review, but an overview of the status of proteome projects has not yet been presented. Advances in proteome projects will initially rely on progress in genome sequencing initiatives, to enable an identity, amino acid sequence, or function to be assigned to each protein spot. Table 7 shows genome size, proteome size, and the number of proteins already defined for a number of model organisms. This indicates that whilst genome sequencing programs for *E. coli* and *S. cerevisiae* are advanced, the massive size of some other genomes (and especially the human genome) means that their complete nucleotide sequences are unlikely to be available for many years. Because of this, 2-D reference maps and proteome projects of single cell organisms like *Mycoplasma sp.*, *E. coli* and *S. cerevisiae* will be the most detailed (Cordwell *et al.*, 1995; Wasinger *et al.*, 1995; Vanbogelen *et al.*, 1992; Garrels *et al.*, 1994), and complete maps of other organisms will take longer to construct. However, the use of cross-species protein identification techniques will allow proteomes of many prokaryotes and simple eukaryotes to be partially defined in reference to *E. coli* and *S. cerevisiae*.

Table 7: Estimated genome size, estimated proteome size, number of protein sequences in SWISS-PROT Release 31 (March, 1995), and approximate number of proteins of known identity on 2-D reference maps for some model organisms. Genome size data from Smith (1994), and total protein data from Bird (1995). Genome sequencing projects of *E. coli* and *S. cerevisiae* will probably be complete in 1996

Species Name	Haploid genome size (million bp)	Estimated proteome size (total proteins)	Protein entries in SWISS-PROT	Proteins annotated on 2-D Maps
<i>Mycoplasma</i> species	0.6–0.8	300–600	100	> 100
<i>Escherichia coli</i>	4.8	4000	3170	> 300
<i>Saccharomyces cerevisiae</i>	12.5	6000	3160	> 100
<i>Drosophila melanogaster</i>	70	12500	2102	–
<i>Arabidopsis thaliana</i>	70	14000	270	–
<i>Caenorhabditis elegans</i>	80	17000	703	–
<i>Homo sapiens</i>	2900	60000–80000	3326	> 1000

The study of vertebrate proteomes and vertebrate development is a phenomenal undertaking in comparison to the investigation of single cell organisms. This is because vast numbers of proteins are developmentally expressed, each body tissue has hundreds of unique proteins, and there are numerous tissue types. However, it is estimated that at least 35% of proteins in vertebrate cells will be conserved from tissue to tissue, constituting the 'housekeeping' proteins (Bird, 1995), with the remainder of proteins constituting a set that are specific to a cell type. Providing that standardised electrophoretic conditions are used, reference maps from many tissues of one organism can be superimposed in gel databases (e.g. Hochstrasser *et al.*, 1992). This accelerates the definition of the 'housekeeping' proteins, as well as sets of proteins that are unique to different tissue types. Such studies may, however, be complicated by post-translational modifications, which can differ on the same gene product in different tissues. Proteins that remain unknown after identification procedures will be useful in providing focus for nucleic acid sequencing initiatives.

## FUTURE DIRECTIONS OF PROTEOME PROJECTS

This review has described recent advances in the area of proteome research. It has illustrated how new developments of older techniques (2-Delectrophoresis, and amino acid analysis) as well as the applications of new technology (mass spectrometry) have greatly widened the choice of tools the biologist and protein chemist has for the separation, identification and analysis of complex mixtures of proteins. This has made possible the establishment of detailed reference maps for organisms, which are becoming the method of choice for the definition of tissues or whole cells, and the investigation of gene expression therein.

Proteome projects are already impacting on the dogma of molecular biology that DNA sequence constitutes the definition of an organism. For example, the proteomes of different tissues of a single organism are often significantly different. Similarly, cross-species identification of proteins (for example the identification of proteins from *Candida albicans* by comparison with *S. cerevisiae*) can open up studies on organisms that are poorly molecularly defined. As cross-species identification can proceed at a pace orders of magnitude faster than a genome project in terms of defining the gene and protein complement of organisms, the need for the DNA sequencing of genomes will be avoided, and emphasis placed on those found to be novel.

Just as genome sequencing is not an end in itself, neither is an annotated 2-D protein reference map of an organism, nor indeed the identification of proteins in a proteome. So whilst an immediate aim of proteome projects is to screen proteins in reference maps, this will lead to expression studies and characterisation of post-translational modifications. The challenge that then needs to be addressed is the investigation of structure and function of proteins in a proteome. The magnitude of this is illustrated by the fact that over half the open reading frames identified in *S. cerevisiae* chromosome III were initially of no known function (Oliver *et al.*, 1992). Structural and functional studies will be an undertaking just as formidable as genome studies are now and proteome projects are becoming, but will lead to an unimaginably detailed understanding of how living organisms are constructed and how they operate.

## Acknowledgements

MRW is the recipient of an Australian Postgraduate Research Award. AAG, MRW, IHS and KIW acknowledge assistance for proteome projects through Macquarie University Research Grants, the Australian Research Council, the Australian National Health and Medical Research Council, Beckman Instruments and GBC Scientific Equipment. DH acknowledges the financial support of a Montux Foundation Grant and the Swiss National Fund for Scientific Research (Grant # 31-33658.92). We thank colleagues who supplied work that was 'In Press' during the writing of this review.

## References

- ANDERSON, N.L., HOFMANN, J.P., GEMMELL, A. AND TAYLOR, J. (1984). Global approaches to quantitative analysis of gene-expression patterns observed by use of two-dimensional gel electrophoresis. *Clinical Chemistry*, **30**, 2031-2036.

- ANDERSON, N.L. AND ANDERSON, N.G. (1991). A two-dimensional gel database of human plasma proteins. *Electrophoresis*, **12**, 853-906.
- ANDERSON, N.L., ESOLER-BLASCO, R., HOFMANN, J.P. AND ANDERSON, N.G. (1991). A two-dimensional gel database of rat liver proteins useful in gene regulation and drug effects studies. *Electrophoresis*, **12**, 907-930.
- ANDERSON, N.L., COPPLE, D.C., BENDELE, R.A., PROBST, G.S., RICHARDSON, F.C. (1992). Covalent protein modifications and gene expression changes in rodent liver following administration of methapyrilene: a study using two-dimensional electrophoresis. *Fundamental and Applied Toxicology*, **18**, 570-580.
- APPEL, R.D., BAIRDOCH, A. AND HOCHSTRASSER, D.F. (1994). A new generation of information retrieval tools for biologists: the example of the EXPASY WWW server. *Trends in Biochemical Sciences*, **19**, 256-260.
- APPEL, R.D., HOCHSTRASSER, D.F., FUNK, M., VARGAS, J.R., PELIGRINI, C., MÜLLER, A.F. AND SCHERRER, J.R. (1991). The MELANIE project: from a biopsy to automatic protein map interpretation by computer. *Electrophoresis*, **12**, 722-735.
- APPEL, R.D., SANCHEZ, J.-C., BAIRDOCH, A., GOLAZ, O., MIL, M., VARGAS, J.R. AND HOCHSTRASSER, D.F. (1993). SWISS-2DPAGE: a database of two-dimensional gel electrophoresis images. *Electrophoresis*, **14**, 1323-1328.
- APPEL, R.D., SANCHEZ, J.-C., BAIRDOCH, A., GOLAZ, O., RAVIER, F., PASQUALI, C., HUGHES, G. AND HOCHSTRASSER, D.F. (1994). The SWISS-2DPAGE database of two-dimensional polyacrylamide gel electrophoresis. *Nucleic Acids Research*, **22**, 3581-3582.
- BAIRDOCH, A. AND BOECKMANN, B. (1994). The SWISS-PROT protein sequence database: current status. *Nucleic Acids Research*, **22**, 3576-3580.
- BAKER, C.S., CORBETT, J.M., MAY, A.J., YACOB, M.H. AND DUNN, M.J. (1992). A human myocardial two-dimensional electrophoresis database: protein characterisation by microsequencing and immunoblotting. *Electrophoresis*, **13**, 723-726.
- BARTLET-JONES, M., JEFFERY, W.A., HANSEN, H.F. AND PAPPIN, D.J.C. (1994). Peptide ladder sequencing by mass spectrometry using a novel, volatile degradation reagent. *Rapid Communications in Mass Spectrometry*, **8**, 737-742.
- BAUER, D., MÜLLER, H., REICH, J., RIEDEL, H., AHRENKIEL, V., WARTHOF, P. AND STRAUSS, M. (1993). Identification of differentially expressed mRNA species by an improved display technique (DDRT-PCR). *Nucleic Acids Research*, **21**, 2272-2280.
- BERNERS-LEE, T.J., CAILLIAT, R., GROFF, J.F. AND POLLERMAN, B. (1992). *Electronic Networks: Research, Applications, and Policy*, **2**, 52-58.
- BIRD, A.P. (1995). Gene number, noise reduction and biological complexity. *Trends in Genetics*, **11**, 9-100.
- BJELLOVIST, B., EK, K., RIGHETTI, P.G., GIANAZZA, E., GORG, A., WESTERMEIER, R. AND POSTEL, W. (1982). Isoelectric focusing in immobilized pH gradients: principle, methodology and some applications. *Journal of Biochemical and Biophysical Methods*, **6**, 317-339.
- BJELLOVIST, B., PASQUALI, C., RAVIER, F., SANCHEZ, J.-C. AND HOCHSTRASSER, D.F. (1993a). A nonlinear wide-range immobilized pH gradient for two-dimensional electrophoresis and its definition in a relevant pH scale. *Electrophoresis*, **14**, 1357-1365.
- BJELLOVIST, B., SANCHEZ, J.-C., PASQUALI, C., RAVIER, F., PAQUET, N., FRITIGER, S., HUGHES, G.J. AND HOCHSTRASSER, D.F. (1993b). Micropreparative 2-D electrophoresis allowing the separation of milligram amounts of proteins. *Electrophoresis*, **14**, 1375-1378.
- BJELLOVIST, B., HUGHES, G., PASQUALI, C., PAQUET, N., RAVIER, F., SANCHEZ, J.-C., FRITIGER, S. AND HOCHSTRASSER, D. (1993c). The focusing positions of polypeptides in immobilized pH gradients can be predicted from their amino acid sequences. *Electrophoresis*, **14**, 1025-1031.
- BOYNER, W.M. AND LASKEY, R.A. (1974). A film detection method for tritium-labeled proteins and nucleic acids in polyacrylamide gels. *European Journal of Biochemistry*, **46**, 83-88.
- BOITTELL, T., GARRELS, J.I., FRANZA, B.R., MONARDO, P.J. AND LATTER, G.I. (1994). REF52: a global gel navigator: an internet-accessible two-dimensional gel electrophoresis database. *Electrophoresis*, **15**, 1487-1490.
- BREWER, J., GRUND, E., HAGERLID, P., OLSSON, I. AND LIZANA, J. (1986). In *Electrophoresis '86* (M.J. Dunn, Ed.), pp. 226-229. VCH, Weinheim.





- Mass spectrometric analysis of blotted proteins after gel electrophoretic separation by matrix-assisted laser desorption/ionization. *Electrophoresis*, 13, 664-668.
- ECKERSKORN, C. AND LOTTSPRECH, F. (1993). Structural characterization of blotting membranes and the influence of membrane parameters for electroblotting and subsequent amino acid sequence analysis of proteins. *Electrophoresis*, 14, 531-536.
- EL, K., BJELLOVIST, B.J. AND RIGHETTI, P.G. (1983). Preparative isoelectric focusing in immobilized pH gradients. I. General principle and methodology. *Journal of Biochemical and Biophysical Methods*, 8, 135-155.
- FEL, S.J., CARLSEN, J., MOSE LARSEN, P., JENSEN, L.A., KJELDSEN, K. AND HANSEN, S. (1994). Two-dimensional gel electrophoresis as a tool for molecular cardiology. Proceedings of the International Society for Heart Research (XV European Section Meeting), pp. 9-16.
- FREY, J.R., KUHN, L., KETTMAN, J.R. AND LEFKOVITS, I. (1994). The amino acid composition of 350 lymphocyte proteins. *Molecular Immunology*, 31, 1219-1231.
- GARRELS, J.I. (1989). The QUEST system for quantitative analysis of two-dimensional gels. *Journal of Biological Chemistry*, 264, 5269-5282.
- GARRELS, J.I. AND FRANZA, B.R. (1989). The REF52 protein database. *Journal of Biological Chemistry*, 264, 5283-5295.
- GARRELS, J.I., FRANZA, B.R., CHANG, C. AND LATTER, G. (1990). Quantitative exploration of the REF52 protein database: cluster analysis reveals the major protein expression profiles in response to growth regulation, serum stimulation, and viral transformation. *Electrophoresis*, 11, 1114-1130.
- GARRELS, J.I., FUTCHER, B., KORIYASHI, R., LATTER, I., SCHWENDER, B., VOLPE, T., WARNER, J.R. AND MCLAUGHLIN, C.S. (1994). Protein identification for a *Saccharomyces cerevisiae* protein database. *Electrophoresis*, 15, 1466-1486.
- GELFI, C., BOSSI, M.L., BJELLOVIST, B. AND RIGHETTI, P.G. (1987). Isoelectric focusing in immobilized pH gradients in the pH 10-11 range. *Journal of Biochemical and Biophysical Research Methods*, 15, 41-48.
- GHARAHDAHLI, F., ATHERTON, D., DEMOTT, M. AND MISCHE, S.M. (1992). Amino acid analysis of PVDF-bound proteins. In *Techniques in Protein Chemistry III* (R.H. Aegele, Ed.), pp. 249-260. Academic Press, San Diego.
- GHOSH, P., OKOH, C., LIU, Q.H. AND LAKSHMAN, M.R. (1993). Effects of chronic ethanol on enzyme-regulating n-acylation and desaturation of transferrin in rats. *Alcoholism: Clinical and Experimental Research*, 17, 576-579.
- GIOMETTI, C.S., GENNIELL, M.A., TOLLAKSEN, S.L. AND TAYLOR, J. (1991). Quantitation of human leukocyte proteins after silver staining: a study with two-dimensional electrophoresis. *Electrophoresis*, 12, 536-543.
- GIOMETTI, C.S., TAYLOR, J. AND TOLLAKSEN, S.L. (1992). Mouse liver protein database: a catalog of proteins detected by two-dimensional gel electrophoresis. *Electrophoresis*, 13, 970-990.
- GOLAZ, O., HUGHES, G.J., FRITIGER, S., PAQUET, N., BAIRICH, A., PASOLI, C., SANCHEZ, J.C., TISSOT, J.D., APPEL, R.D., WALZER, C., BALANT, L. AND HOCHSTRASSER, D.F. (1993). Plasma and red blood cell protein maps: update 1993. *Electrophoresis*, 14, 1223-1231.
- GOLD, M.R., YUNGWIRTH, T., SUTHERLAND, C.L., INGHAM, R.J., VIANZON, D., CHIL, R., VAN OOSTVEEN, L., MORRISON, H.D. AND AEBERSOLD, R. (1994). Purification and identification of tyrosine-phosphorylated proteins from lymphocytes stimulated through the antigen receptor. *Electrophoresis*, 15, 441-453.
- GOLDHARG, H.A., DOMENICUCCI, C., PRINGLE, G.A. AND SODEK, J. (1988). Mineral-binding proteoglycans of fetal porcine calvarial bone. *Journal of Biological Chemistry*, 263, 12092-12101.
- GOOLEY, A.A., MARSHCHALEK, R. AND WILLIAMS, K.L. (1992). Size polymorphism due to changes in the number of O-glycosylated tandem repeats in the *Drosophila discaloid* glycoprotein P-A. *Genetics*, 130, 749-756.
- GORG, A., POSTEL, W. AND GUNTHER, S. (1988). The current state of two-dimensional electrophoresis with immobilized pH gradients. *Electrophoresis*, 9, 531-546.
- GORG, A., POSTEL, W., GUNTHER, S., WESER, J., STRÄHLER, J.R., HANASHI, S.M., SOMERLOT, L.

- AND KLICK, R. (1988). Approach to stationary two-dimensional pattern influence of focusing time and iminobulin carrier ampholyte concentrations. *Electrophoresis*, 9, 37-46.
- GRAVEL, P., GOLAZ, O., WALZER, C., HOCHSTRASSER, D.F., TURLEZ, H., AND BALANT, L.P. (1994). Analysis of glycoproteins separated by two-dimensional gel electrophoresis using lectin blotting revealed by chemiluminescence. *Analytical Biochemistry*, 221, 66-71.
- GUTHRIE, S., POSTEL, W., WILKING, H., AND GORG, A. (1988). Acid phosphatase typing for breeding nematode-resistant tomatoes by isoelectric focusing with an ultranarrow immobilized pH gradient. *Electrophoresis*, 9, 618-620.
- HANASH, S.M., STRAHLER, J.R., NEEL, J.V., HAILAT, N., MELHEM, R., KEIM, D., ZHANG, X.X., WAGNER, D., GAGE, D.A. AND WATSON, J.T. (1991). Highly resolving two-dimensional gels for protein sequencing. *Proceedings of the National Academy of Sciences USA*, 88, 5709-5713.
- HARRINGTON, M.G., COFFMAN, J.A., CALZONE, F.J., HOOD, L.E., BRITTON, R.J. AND DAVIDSON, E.H. (1992). Complexity of sea urchin embryo nuclear proteins that contain basic domains. *Proceedings of the National Academy of Sciences USA*, 89, 6252-6256.
- HARRINGTON, M.G., LEE, K.H., YUN, M., ZEVERT, T., BAILEY, J.E. AND HOOD, L.E. (1993). Mechanical precision in two-dimensional electrophoresis can improve spot positional reproducibility. *Applied and Theoretical Electrophoresis*, 3, 347-353.
- HELLMAN, U., WERNSTEDT, C., GONZALEZ, J. AND HELDIN, C.-H. (1995). Improvement of an in-gel digestion for the micropreparation of internal protein fragments for amino acid sequencing. *Analytical Biochemistry*, 224, 451-455.
- HEYNEL, W.J., BILLECI, T.M., STULTS, J.T., WONG, S.C., GRINLEY, C. AND WATANABE, C. (1993). Identifying proteins from two-dimensional gels by molecular mass searching of peptide fragments in protein sequence databases. *Proceedings of the National Academy of Sciences USA*, 90, 5011-5015.
- HESS, D., COVEY, T.C., WINZ, R., BROWNSEY, R.W. AND AEBERSOLD, R. (1993). Analytical and micropreparative peptide mapping by high performance liquid chromatography/electrospray mass spectrometry of proteins purified by gel electrophoresis. *Protein Science*, 2, 1342-1351.
- HOBOMI, U., HOLTHAEVE, T. AND SANDER, C. (1994). Amino acid analysis and protein database compositional search as a rapid and inexpensive method to identify proteins. *Analytical Biochemistry*, 222, 202-209.
- HOCHSTRASSER, D.F. AND MERRIL, C.R. (1988). 'Catalysts' for polyacrylamide gel polymerization and detection of proteins by silver staining. *Applied and Theoretical Electrophoresis*, 1, 35-40.
- HOCHSTRASSER, D.F., PATCHORNIAK, A. AND MERRIL, C.R. (1988). Development of polyacrylamide gels that improve the separation of proteins and their detection by silver staining. *Analytical Biochemistry*, 173, 412-423.
- HOCHSTRASSER, A.C., JAMES, R.W., POMETTA, D. AND HOCHSTRASSER, D.F. (1991a). Preparative isoelectric focusing and high resolution two-dimensional electrophoresis for concentration and purification of proteins. *Applied and Theoretical Electrophoresis*, 1, 333-337.
- HOCHSTRASSER, D.F., APPEL, R.D., VARGAS, R., PERKIER, R., VILLOD, J.F., RAVIER, F., PASQUALI, C., FUNK, M., PELLIGRINI, C., MÜLLER, A.F. AND SCHIERRER, J.R. (1991b). A clinical molecular scanner: the Melanin project. *Medical Computing*, 8, 85-91.
- HOCHSTRASSER, D.F., FRITZGER, S., PAQUET, N., BAIRIOCHI, A., RAVIER, F., PASQUALI, C., SANCHEZ, J.-C., TISSOT, J.-D., BJELLOVIST, B., VARGAS, R., APPEL, R.D. AND HUGHES, G.J. (1992). Human liver protein map: a reference database established by microsequencing and gel comparison. *Electrophoresis*, 13, 992-1001.
- HOLT, T.G., CHANG, C., LAURENT-WINTER, C., MURAKAMI, T., DAVIES, J.E. AND THOMPSON, C.J. (1992). Global changes in gene expression related to antibiotic synthesis in *Streptomyces hygroscopicus*. *Molecular Microbiology*, 6, 969-980.
- HONORE, B., LEFFERS, H., MADSEN, P. AND CELIS, J.E. (1993). Interferon-gamma up-regulates a unique set of proteins in human keratinocytes. Molecular cloning and expression of the cDNA encoding the RGD-sequence containing protein IGUP1-5111. *European Journal of Biochemistry*, 218, 421-430.
- HILBERT, M.C., VATH, J.E., YU, W. AND MARTIN, S.A. (1993). Site-specific carbohydrate

- identification in recombinant proteins using MALD-TOF MS. *Analytical Chemistry*, 65, 2791-2800.
- HUGHES, G.J., FRUTIGER, S., PAQUET, N., PASOLUNghi, C., SANCHEZ, J.-C., TISSOT, J.D., BAIRIOCH, A., APPEL, R.D. AND HOCHSTRASSER, D.F. (1993). Human liver protein map update 1993. *Electrophoresis*, 14, 1216-1222.
- HUGHES, J.H., MACK, K. AND HAMPARIAN, V.V. (1988). India ink staining of proteins on nylon and hydrophobic membranes. *Analytical Biochemistry*, 173, 18-25.
- JAMES, P., QI, ADRONI, M., CARAFOLI, E. AND GONNET, G. (1993). Protein identification by mass profile fingerprinting. *Biochemical and Biophysical Research Communications*, 195, 56-64.
- JIL, H., WHITEHEAD, R.H., REID, G.E., MORITZ, F.L., WARD, L.D. AND SIMPSON, R.J. (1994). Two-dimensional electrophoretic analysis of proteins expressed by normal and cancerous human crypts: application of mass spectrometry to peptide-mass fingerprinting. *Electrophoresis*, 15, 391-405.
- JOHNSON, R.S. AND WALSH, K.A. (1992). Sequence analysis of peptide mixtures by automated integration of Edman and mass spectrometric data. *Protein Science*, 1, 1083-1091.
- JOHNSTON, R.F., PICKETT, S.C. AND BARKER, D.L. (1990). Autoradiography using storage phosphor technology. *Electrophoresis*, 11, 355-360.
- JUNGILLT, P., DZIONARA, M., KLOSE, J. AND WITTMANN-LEIBOLD, B. (1992). Identification of tissue proteins by amino acid analysis after purification by two-dimensional electrophoresis. *Journal of Protein Chemistry*, 11, 603-612.
- JUNGILLT, P., OTTO, A., ZEINDL-EBERHART, E., PLEIBNER, K.-P., KNECHT, M., REGITZ-ZAGROSEK, V., FLECK, E. AND WITTMANN-LEIBOLD, B. (1994). Protein composition of the human heart: the construction of a myocardial two-dimensional electrophoresis database. *Electrophoresis*, 15, 685-707.
- KOHARA, Y., AKIYAMA, K. AND ISONO, K. (1987). The physical map of the whole *E. coli* chromosome: application of a new strategy for rapid analysis and sorting of a large genomic library. *Cell*, 50, 495-508.
- KLOSE, J. (1975). Protein mapping by combined isoelectric focusing and electrophoresis in mouse tissues. A novel approach to testing for individual point mutations in mammals. *Human Genetics*, 26, 231-243.
- LATHAM, K.E., GARRELS, J.I., CHANG, C. AND SOLTER, D. (1991). Quantitative analysis of protein synthesis in mouse embryos I: extensive re-programming at the one- and two-cell stages. *Development*, 2, 921-932.
- LATHAM, K.E., GARRELS, J.I., CHANG, C. AND SOLTER, D. (1992). Analysis of embryonic mouse development: construction of a high-resolution, two-dimensional gel protein database. *Applied and Theoretical Electrophoresis*, 2, 163-170.
- LATHAM, K.E., GARRELS, J.I. AND SOLTER, D. (1993). Two-dimensional analysis of protein synthesis. *Methods in Enzymology*, 225, 473-489.
- LE MAIRE, M., DESCHAMPS, S., MOLLER, J.V., LE CAER, J.P. AND ROUSIER, J. (1993). Electrospray ionization mass spectrometry from sodium dodecyl sulfate-polyacrylamide gel electrophoresis: application to the topology of the sarcoplasmic reticulum  $\text{Ca}^{2+}$  ATPase. *Analytical Biochemistry*, 214, 50-57.
- LEVIN, P.F. AND LESTER, E.P. (1989). Database and search techniques for two-dimensional gel protein data: a comparison of paradigms for exploratory data analysis and prospects for biological modelling. *Electrophoresis*, 10, 122-140.
- LEVIN, P.F., WU, Y. AND UPTON, K. (1993). An efficient disk-based data structure for rapid searching of quantitative two-dimensional gel databases. *Electrophoresis*, 14, 1341-1350.
- LI, K.W., GERAERTS, W.P., VAN-ELK, R. AND KOOSE, J. (1989). Quantification of proteins in the subnanogram and nanogram range: comparison of the AutoDye, FerroDye, and India ink staining methods. *Analytical Biochemistry*, 182, 33-47.
- LIANG, P. AND PARDEE, A.B. (1992). Differential display of eukaryotic messenger RNA by means of the polymerase chain reaction. *Science*, 257, 967-971.
- MANN, M. (1993). Sequence database searching by mass spectrometric data. In *Microcharacterisation of Proteins* (R. Kellner, F. Lottspeich, and H.E. Meyer, Eds.), pp 223-245. VCH, Weinheim.

- MANN, M., HOJRLU, P. AND ROEPSTORFF, P. (1993). Use of mass spectrometric molecular weights to identify proteins in sequence databases. *Biological Mass Spectrometry*, **22**, 335-345.
- MANN, M. AND WILKINS, M. (1994). Error tolerant identification of peptides in sequence databases by peptide sequence tags. *Analytical Biochemistry*, **266**, 4390-4399.
- MATSUDAIRA, P. (1987). Sequence of picomole quantities of proteins electroblotted onto polyvinylidene difluoride membranes. *Journal of Biological Chemistry*, **262**, 10035-10038.
- MONARDO, P.J., BOUTELL, T., GARRELS, J.I. AND LATTER, G.I. (1994). A distributed system for two-dimensional gel analysis. *Computer Applications in the Biosciences*, **10**, 137-143.
- MORTZ, E., VOPM, O., MANN, M. AND ROEPSTORFF, P. (1994). Identification of proteins in polyacrylamide gels by mass spectrometric peptide mapping combined with database search. *Biological Mass Spectrometry*, **23**, 249-261.
- MURPHY, L.R. AND LOBAL, K. (1991). Measurement of picomoles of phosphoamino acids by high performance liquid chromatography. *Analytical Biochemistry*, **193**, 299-303.
- MYRICK, J.E., LENKIN, P.F., ROBINSON, M.K. AND UPTON, K.M. (1993). Comparison of the Biomeq Vantage 2000 and the GELLAB-II two-dimensional electrophoresis image analysis systems. *Applied and Theoretical Electrophoresis*, **3**, 333-346.
- NEIDHARDT, F.C., APPLEBY, D.B., SANKAR, P., HUTTON, M.E. AND PHILLIPS, T.A. (1989). Genomically linked cellular protein databases derived from two-dimensional polyacrylamide gel electrophoresis. *Electrophoresis*, **10**, 116-122.
- NIKODEM, V. AND FRESCO, J.R. (1979). Protein fingerprinting by SDS-gel electrophoresis after partial fragmentation with CNBr. *Analytical Biochemistry*, **97**, 382-386.
- NOKIHARA, K., MORITA, N. AND KURIKI, T. (1992). Applications of an automated apparatus for two-dimensional electrophoresis. Model TEP-1, for microsequence analysis of proteins. *Electrophoresis*, **13**, 701-707.
- O'FARRELL, P.H. (1975). High resolution two-dimensional electrophoresis of proteins. *Journal of Biological Chemistry*, **250**, 4007-4021.
- O'FARRELL, P.Z., GOODMAN, H.M. AND O'FARRELL, P.H. (1977). High resolution two-dimensional electrophoresis of basic as well as acidic proteins. *Cell*, **12**, 1133-1142.
- OLIVER, A.L. (1992). The complete DNA sequence of yeast chromosome III. *Nature*, **357**, 38-46.
- OLSEN, A.D. AND MILLER, M.J. (1988). Elsie 4: quantitative computer analysis of sets of two-dimensional gel electrophoretograms. *Analytical Biochemistry*, **169**, 49-70.
- ORTIZ, M.L., CALERO, M., FERNANDEZ-PATRON, C., PATRON, C.F., CASTELLANOS, L. AND MENDEZ, E. (1992). Imidazole-SDS-Zn reverse staining of proteins in gels containing or not SDS and microsequence of individual unmodified electroblotted proteins. *FEBS Letters*, **296**, 300-304.
- OSTERGREN, K., ERIKSSON, G. AND BJELLOVIST, B. (1988). The influence of support material used on band sharpness in Immobiline gels. *Journal of Biochemical and Biophysical Methods*, **16**, 165-170.
- OL, K., WILKINS, M.R., YAN, J.X., GOOLEY, A.A., FUNG, Y., SHELMACK, D. AND WILLIAMS, K.L. (1995). Improved high-performance liquid chromatography of amino acids derivatized with 9-fluorenylmethyl chloroformate. *Journal of Chromatography* (in press).
- PACKER, N., WILKINS, M.R., GOLAZ, O., LAWSON, M., GOOLEY, A.A., HOCHSTRASSER, D.F., REDMOND, J. AND WILLIAMS, K.L. (1995). Characterisation of human plasma glycoproteins separated by two-dimensional gel electrophoresis. *BioTechnology* (in press).
- PAPPIN, D.J.C., HOJRLU, P. AND BLEASBY, A.J. (1993). Rapid identification of proteins by peptide-mass fingerprinting. *Current Biology*, **3**, 327-332.
- PATTERSON, S.D. (1994). From electrophoretically separated protein to identification: strategies for sequence and mass analysis. *Analytical Biochemistry*, **221**, 1-15.
- PATTERSON, S.D. AND LATTER, G.I. (1993). Evaluation of storage phosphor imaging for quantitative analysis of 2-D gels using the Quest II system. *BioTechniques*, **15**, 1076-1083.
- PISANO, A., REDMOND, J.W., WILLIAMS, K.L. AND GOOLEY, A.A. (1993). Glycosylation sites identified by solid-phase Edman degradation: O-linked glycosylation motifs on human glycoprotein A. *Glycobiology*, **3**, 429-435.
- RABILLOI, D.T. (1992). A comparison between low background silver diamine and silver nitrate protein stains. *Electrophoresis*, **13**, 429-439.

- RASMUSSEN, H.H., VAN DAMME, J., PUYE, M., GESSER, B., CELIS, J.E. AND VANDEKERCKHOVE, J. (1992). Microsequences of 145 proteins recorded in the two-dimensional gel protein database of normal human epidermal keratinocytes. *Electrophoresis*, 13, 960-969.
- RASMUSSEN, H.H., MORTZ, E., MANN, M., ROEPSTORFF, P. AND CELIS, J.E. (1994). Identification of transformation sensitive proteins recorded in human two-dimensional gel protein databases by mass-spectrometric peptide mapping alone and in combination with microsequencing. *Electrophoresis*, 15, 406-416.
- RICHARDSON, F.C., HORN, D.M. AND ANDERSON, N.L. (1994). Dose-responses in rat hepatic protein modification and expression following exposure to the rat hepatocarcinogen methapyrilene. *Carcinogenesis*, 15, 325-329.
- RIGHETTI, P.G. (1990). Immobilized pH gradients: theory and methodology. In *Laboratory Techniques in Biochemistry and Molecular Biology* (R.H. Burdon and P.H. van Knippenberg, Eds) Elsevier, Amsterdam.
- RIGHETTI, P.G. AND DRYSDALE, J.W. (1973). *Annals of the New York Academy of Sciences*, 209, 163-186.
- RODRIGUEZ, L.V., GERNSTEN, D.M., RAMAGLI, L.S. AND JOHNSTON, D.A. (1993). Towards stoichiometric silver staining of proteins resolved in complex two-dimensional electrophoresis gels: real-time analysis of pattern development. *Electrophoresis*, 14, 628-637.
- ROSENFELD, J., CAPDEVILLE, J., GUILLENOT, J.C. AND FERRARA, P. (1992). In-gel digestion of proteins for internal sequence analysis after one- or two-dimensional gel electrophoresis. *Analytical Biochemistry*, 203, 173-179.
- SANCHEZ, J.C., RAVIER, F., PASOLUNGI, C., FRUTIGER, S., PAQUET, N., BJELLOVIST, B., HOCHSTRASSER, D.F. AND HUGHES, G.J. (1992). Improving the detection of proteins after transfer to polyvinylidene difluoride membranes. *Electrophoresis*, 13, 715-717.
- SANGER, F., COLLISON, A.R., HONG, G.F., HILL, D.F. AND PETERSEN, G.B. (1962). Nucleotide sequence of bacteriophage  $\lambda$  DNA. *Journal of Molecular Biology*, 162, 729-773.
- SCHHEEL, G.J. (1975). Two-dimensional analysis of soluble proteins. *Biochemistry*, 250, 5375-5385.
- SHAW, G. (1993). Rapid identification of proteins. *Proceedings of the National Academy of Sciences USA*, 90, 5138-5142.
- SIBBALD, P.R., SOMMERFELDT, H. AND ARGOS, P. (1991). Identification of proteins in sequence databases from amino acid composition. *Analytical Biochemistry*, 198, 330-333.
- SIMPSON, R.J., TSUGITA, A., CELIS, J.E., GARRELS, J.I. AND MEWES, H.W. (1992). Workshop on two-dimensional gel protein databases. *Electrophoresis*, 13, 1055-1061.
- SINHA, P.K., KOTTGEV, E., STOFFLER, M.M., GIANAZZA, E. AND RIGHETTI, P.G. (1990). Two-dimensional maps in very acidic immobilized pH gradients. *Journal of Biochemical and Biophysical Methods*, 20, 345-352.
- SMITH, D.W. (1994). Introduction. In *Bioinformatics: Informatics and Genome Projects* (D.W. Smith, Ed.), pp1-12. Academic Press, San Diego.
- STRUPAT, K., KARAS, M., HILLENKAMP, F., ECKERSKORN, C. AND LUTTSPEICH, F. (1994). Matrix-assisted laser desorption/ionization mass spectrometry of proteins electrophoresed after polyacrylamide gel electrophoresis. *Analytical Chemistry*, 66, 466-470.
- SUTTON, C.W., PENIBERTON, K.S., COTTELL, J.S., CORBETT, J.M., WHEELER, C.H., DUNN, M.J. AND PAPPIN, D.J. (1995). Identification of myocardial proteins from two-dimensional gels by peptide mass fingerprinting. *Electrophoresis*, 16, 306-310.
- TOLS, G.J., FAUSNAUGH, J.L., AKINYOYE, O., LACKLAND, H., WINTERCASH, P., VITORICA, F.J. AND STEIN, S. (1989). Amino acid analysis on polyvinylidene difluoride membranes. *Analytical Biochemistry*, 179, 50-55.
- TOVEY, E.R., FORD, S.A. AND BALDO, B.A. (1987). Protein blotting on nitrocellulose: some important aspects of the resolution and detection of antigens in complex extracts. *Journal of Biochemical and Biophysical Methods*, 14, 1-17.
- URWIN, V.E. AND JACKSON, P. (1993). Two-dimensional polyacrylamide gel electrophoresis of proteins labeled with the fluorophore monobromobimane prior to first-dimensional isoelectric focusing: imaging of the fluorescent protein spot patterns using a cooled charge-coupled device. *Analytical Biochemistry*, 209, 57-62.
- VAN BOGHELEN, R.A., HUTTON, M.E. AND NEIDHARDT, F.C. (1990). Gene-protein database

- of *Escherichia coli*. K-12, edition 3. *Electrophoresis*, 11, 1131-1166.
- VANBOGELEN, R.A. AND NEIDHARDT, F.C. (1991). The gene-protein database of *Escherichia coli*, edition 4. *Electrophoresis*, 12, 955-994.
- VANBOGELEN, R.A., SANKER, F., CLARK, R.L., BOGAN, J.A. AND NEIDHARDT, F.C. (1992). The gene-protein database of *Escherichia coli*, edition 5. *Electrophoresis*, 13, 1011-1054.
- VANDEKERKHOVE, J., BAUW, G., VANCOPIERNOLE, K., HONORE, B. AND CELIS, J. (1990). Comparative two-dimensional gel analysis and microsequencing identifies gel-solin as one of the most prominent downregulated markers of transformed human fibroblast and epithelial cells. *Journal of Cell Biology*, 111, 95-102.
- VANPLETEREN, J.R., RAYMACKERS, J.G., VAN BUN, S.M. AND MEHUS, L.A. (1992). Peptide mapping and microsequencing of proteins separated by SDS-PAGE after limited *in situ* hydrolysis. *BioTechniques*, 12, 550-557.
- VORM, O. AND MANN, M. (1994). Improved mass accuracy in matrix-assisted laser desorption/ionization time-of-flight mass spectrometry of peptides. *Journal of the American Society for Mass Spectrometry*, 5, 955-958.
- VORM, O., ROEPSTORFF, P. AND MANN, M. (1994). Improved resolution and very high sensitivity in MALDI TOF of matrix surfaces made by fast evaporation. *Analytical Chemistry*, 66, 3281-3287.
- WALLACE, A. AND SALLZ, H.P. (1992a). Ultramicrodetection of proteins in polyacrylamide gels. *Analytical Biochemistry*, 203, 27-34.
- WALLACE, A. AND SALLZ, H.P. (1992b). Beyond silver staining. *Nature*, 357, 605-609.
- WALSH, B.J., GOOLEY, A.A., WILLIAMS, K.L. AND BRETT, S.N. (1995). Identification of macrophage activation associated proteins by two-dimensional electrophoresis and microsequencing. *Journal of Leukocyte Biology*, 57, 507-512.
- WASINGER, V.C., CORDWELL, S.J., POLJAK, A., YAN, J.X., GOOLEY, A.A., WILKINS, M.R., DUNCAN, M., HARRIS, R., WILLIAMS, K.L. AND HUMPHRY-SMITH, I. (1995). Progress with Gene-Product Mapping of the Molluscs: *Myxolusmus genitalium*. *Electrophoresis*, 16, In Press.
- WEITZHANDLER, M., KADLECEK, D., AYDALOVIC, N., FORTE, J.G., CHOW, D. AND TOWNSEND, R.R. (1993). Monosaccharide and oligosaccharide analysis of proteins transferred to polyvinylidene fluoride membranes after sodium dodecyl sulfate-polyacrylamide gel electrophoresis. *Journal of Biological Chemistry*, 268, 5121-5130.
- WILKINS, M.R., PASOLALI, C., APPEL, R.D., OU, K., GOLAZ, O., SANCHEZ, J.-C., YAN, J.X., GOOLEY, A.A., HUGHES, G., HUMPHRY-SMITH, I., WILLIAMS, K.L. AND HOCHSTRASSER, D.F. (1995). From Proteins to Proteomes: large scale protein identification by two-dimensional electrophoresis and amino acid analysis. Submitted.
- WILKINS, M.R., OU, K., APPEL, R.D., GOLAZ, O., PASOLALI, C., YAN, J.X., FARNSWORTH, V., CARTIER, P., HOCHSTRASSER, D.F., WILLIAMS, K.L. AND GOOLEY, A.A. (1996). Rapid protein identification using N-terminal sequence tagging and amino acid analysis (submitted).
- WIRTH, P.J., LYO, L.D., FUJIMOTO, Y., BISGAARD, H.C. AND OLSEN, A.D. (1991). The rat liver epithelial (RLE) cell protein database. *Electrophoresis*, 12, 931-954.
- WIRTH, P.J., LYO, L.D., BENJAMIN, T., HUANG, T.N., OLSEN, A.D. AND PARMALIE, D.C. (1993). The rat liver epithelial (RLE) cell nuclear protein database. *Electrophoresis*, 14, 1199-1215.
- WU, Y., LEMKIN, P.F. AND UPTON, K. (1993). A fast spot segmentation algorithm for two-dimensional gel electrophoresis analysis. *Electrophoresis*, 14, 1351-1356.
- YAMAGUCHI, K. AND ASAKAWA, H. (1988). Preparation of colloidal gold for staining proteins electrotransferred onto nitrocellulose membranes. *Analytical Biochemistry*, 172, 104-107.
- YAMASHITA, K., IDEO, H., OHKURA, T., FUKUSHIMA, K., YASUDA, I., OHNO, K. AND TAKESHITA, K. (1993). Sugar chains of serum transferrin from patients with carbohydrate deficient glycoprotein syndrome. Evidence of asparagine-N-linked oligosaccharide transfer deficiency. *Journal of Biological Chemistry*, 268, 5783-5789.
- YATES, J.R. III, SPEICHER, S., GRIFFIN, P.R. AND HUNKAPILLER, T. (1993). Peptide mass maps: a highly informative approach to protein identification. *Analytical Biochemistry*, 214, 397-408.

5  
The

R.M.F.

Therm  
Hamil  
Wales  
Private

Introduct  
Protein  
dealing  
type. T  
structur  
since th  
by War  
1930) c  
possibly  
mechar  
subtilis  
1986: I  
comme  
proteas  
the cha  
been in  
regulat  
appare  
Kalisz  
Barrett  
Give  
surpris  
attract  
subject  
stabilit  
enzym  
Sherod

\* Contents  
Biochem  
1264-872

## Human cellular protein patterns and their link to genome DNA sequence data: usefulness of two-dimensional gel electrophoresis and microsequencing

JULIO E. CELIS,<sup>1</sup>\* HANNE H. RASMUSSEN,<sup>1</sup> HENRIK LEFFERS,<sup>1</sup> PEDER MADSEN,<sup>1</sup> BENT HONORÉ,<sup>1</sup> BORBALA GESSER,<sup>1</sup> KURT DEJGAARD,<sup>2</sup> JOEL VANDEKERCKHOVE<sup>1</sup>

<sup>1</sup>Institute of Medical Biochemistry and Human Genome Research Centre, Aarhus University, DK-8000 Århus, Denmark and <sup>2</sup>Laboratorium voor Fysiologische Chemie, Rijksuniversiteit Gent, Belgium

**ABSTRACT** Analysis of cellular protein patterns by computer-aided 2-dimensional gel electrophoresis together with recent advances in protein sequence analysis have made possible the establishment of comprehensive 2-dimensional gel protein databases that may link protein and DNA information and that offer a global approach to the study of the cell. Using the integrated approach offered by 2-dimensional gel protein databases it is now possible to reveal phenotype specific protein (or proteins), to microsequence them, to search for homology with previously identified proteins, to clone the cDNAs, to assign partial protein sequence to genes for which the full DNA sequence and the chromosome location is known, and to study the regulatory properties and function of groups of proteins that are coordinately expressed in a given biological process. Human 2-dimensional gel protein databases are becoming increasingly important in view of the concerted effort to map and sequence the entire genome. — Celis, J. E.; Rasmussen, H. H.; Leffers, H.; Madsen, P.; Honoré, B.; Gesser, B.; Dejgaard, K.; Vandekerckhove, J. Human cellular protein patterns and their link to genome DNA sequence data: usefulness of two-dimensional gel electrophoresis and microsequencing. *FASEB J.* 5: 2200-2208; 1991.

**Key Words:** human protein patterns • 2-dimensional gel protein databases • gene expression • microsequencing • cDNA cloning • linking protein and DNA information • genome mapping and sequencing

PROTEINS SYNTHESIZED FROM information contained in the DNA orchestrate most cellular functions. The total number of proteins synthesized by a typical human cell is unknown although current estimates range from 3000 to 6000. Of these, as many as 70% may perform household functions and are expected to be shared by all cell types irrespective of their origin. There are many different cell types in the human body, with perhaps 30,000 to 50,000 proteins expressed in the organism as a whole judged from the fact that about 3% of the haploid genome correspond to genes. Today only a small fraction of the total set of proteins has been identified, and little is known about the protein patterns of individual cell types or their variation under physiological and abnormal conditions.

For the past 15 years, high resolution 2-dimensional gel electrophoresis has been the technique of choice to determine the protein composition of a given cell type and for monitoring changes in gene activity through quantitative and qualitative analysis of the thousands of proteins that orchestrate various cellular functions (refs 1-6 and references

therein). The technique originally described by O'Farrell<sup>1</sup> separates proteins in terms of their isoelectric point (pI) and molecular weight. Usually one chooses a condition of interest and the cell reveals the global protein behavioral response as all detected proteins can be analyzed both qualitatively and quantitatively in relation to each other. At present, most available 2-dimensional gel techniques (regular gel format) can resolve between 1000 and 2000 proteins from a given mammalian cell type, a number that corresponds to about 2 million base pairs of coded DNA. Less abundant proteins can be detected by analyzing partially purified cellular fractions.

Two-dimensional gel electrophoresis has been widely applied to analysis of cellular protein patterns from bacteria to mammalian cells (refs 1-6, and references therein). In spite of much work, however, information gathered from these studies has not reached the scientific community in its fullness because of lack of standardized gel systems and the lack of means for storing and communicating protein information. Only recently, because of the development of appropriate computer software (7-13), has it been possible to store, assign numbers to individual proteins, and store the wealth of information in quantitative and qualitative comprehensive 2-dimensional gel protein databases (4, 14-23), i.e., those containing information about the various properties (physical, chemical, biological, biochemical, physiological, genetic, immunological, architectural, etc.) of all the proteins that can be detected in a given cell type. Such integrated 2-dimensional gel protein databases offer an easy and standardized medium in which to store and communicate protein information and provide a unique framework in which to focus a multidisciplinary approach to study the cell. Once a protein is identified in the database, all of the information accumulated can be easily retrieved and made available to the researcher. In the long run, protein databases are expected to foster a wide variety of biological information that may be instrumental to researchers working in many areas of biology—among others, cancer and oncogene studies, differentiation, development, drug development and testing, genetic variation, and diagnosis of genetic and clinical diseases (Fig. 1).

The approach using systematic 2-dimensional gel protein analysis has recently gained a new dimension with the advent of techniques to microsequence major proteins recorded

\*To whom correspondence should be addressed, at: Institute of Medical Biochemistry and Human Genome Research Centre, Ole Worms Alle, Bldg. 170, University Park, DK-8000 Aarhus C, Denmark.

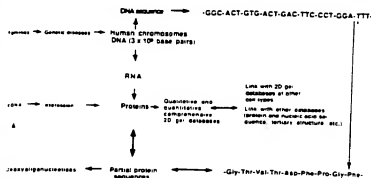


Figure 1. Interface between partial protein sequence databases, comprehensive 2-dimensional gel databases, and the human genome sequencing project. Appropriate software is required to compare protein and DNA sequences. In general, although the inference of a protein's sequence from the DNA sequence (thick arrow) is direct and unambiguous, the DNA sequence can only be inferred approximately from the protein sequence (thin arrow) and cloning of the gene requires either a cDNA or the requisite group of oligonucleotide probes deduced from the partial amino acid sequence. Modified from ref. 6.

in the databases (refs 24-42 and references therein). Partial protein sequences can be used to search for protein identity as well as to prepare specific DNA probes for cloning as yet-uncharacterized proteins (Fig. 1). As these sequences can be stored in the database (see for example Fig. 2H), they offer a unique opportunity to link information on proteins with the existing or forthcoming DNA sequence data on the human genome (Fig. 1) (20, 36, 39).

Using the integrated approach offered by comprehensive 2-dimensional gel databases (Fig. 1), it will be possible to identify phenotype-specific proteins; microsequence them and store the information in the database; search for homology with previously characterized proteins; clone the cDNAs, assign partial protein sequences to genes for which the full DNA sequence and the chromosome location are known, and study the regulatory properties and function of groups of proteins (pathways, organelles, etc.) that are coordinately expressed in a given biological process. Comprehensive 2-dimensional gel protein databases will depict an integrated picture of the expression levels and properties of the thousands of protein components of organelles, pathways, and cytoskeletal systems in both physiological and abnormal conditions and are expected to lead to identification of new regulatory networks in different cell types and organisms. In the future, 2-dimensional gel protein databases may be linked to each other as well as to national and international specialized databanks on nucleic acid and protein sequences, protein structures, NMR experimental data, complex carbohydrates, etc.

A few 2-dimensional gel protein databases that are accessible in a computer form have been published in extenso: these correspond to the protein-gene database of *Escherichia coli* K-12 developed by Neidhardt and colleagues (14, 23), the rat REF 32 database established by Garrels and co-workers at Cold Spring Harbor (18, 22), and a few human databases (transformed amnion cells [15, 20], normal embryonal lung MRC-5 fibroblasts [17, 21], keratinocytes [19] and peripheral blood mononuclear cells [15]) developed in Aarhus. Given space limitations and to keep this review in focus, we will concentrate on the computerized analysis of human cellular 2-dimensional gel patterns, and in particular on the steps involved in establishing comprehensive 2-dimensional gel databases that can link protein and DNA information.

## MAKING AND MANAGING A COMPREHENSIVE 2-DIMENSIONAL GEL DATABASE OF HUMAN CELLULAR PROTEINS

The first step in making a comprehensive 2-dimensional gel protein database is to prepare a synthetic image (digital term of the gel image) of the gel (fluorogram, Coomassie blue or silver stained gel) to be used as a standard or master reference. This can be done with laser scanners, charge couple device (CCD)<sup>2</sup> array scanners, television cameras, rotating drum scanners, and multiwire chambers (13). Computerized analysis systems for spot detection, quantitation, pattern matching, and data handling (access and retrieval of information, database making) have been described in the literature (ELSIE [43], GELLAB [11], HERMES [44], MELANIE [10], QUEST (9), and TYCHO [8]) and some are available commercially (PDQUEST, Protein Database Inc., Huntington, N.Y.; KEPLER, Large Scale Biology, Rockville, Md.; Visage, Biologie Corporation, Ann Arbor, Mich.; Gemini, Joyce Loeb, Gateshead; Microscan 1000, Technology Resources Inc., Nashville, Tenn.; and MasterScan, Billerica, Mass.). Unfortunately, most of these systems are incompatible with one another and their advantages and disadvantages have been discussed by Miller (13).

In our work station in Aarhus, fluorograms are scanned with a Molecular Dynamics laser scanner and the data are analyzed using the PDQUEST II software (Protein Database Inc.) (12) running on a sparc station computer 4100 FC-8-P3 from SUN Microsystems, Inc. The scanner measures intensity in the range of 0-2.0 absorbance. A typical scan of a 17 × 17 cm fluorogram takes about 2 min. Steps in image analysis include: initial smoothing, background subtraction, final smoothing, spot detection, and fitting of ideal Gaussian distribution to spot centers. Spot intensity is calculated as the integration of a fitted Gaussian. If calibration strips containing individual segments of a known amount of radioactivity are used, it is possible to merge multiple exposures of the sample image into a single data image of greater dynamic range. Once the synthetic image is created it can be stored on disk and displayed directly on the monitor. Functions that can be used to edit the images include: cancel (for example, to erase scratches that may have been interpreted as spots by the computer; cancel streaks or low dpm spots), combine (sometimes a spot may be resolved into several closely packed spots), restore, uncombine, and add spot to the gel. The process is time consuming—about 1-1/2 day per image. Edited standard images can be matched to other synthetic images. Figure 2A shows a portion of a standard synthetic image (IEF) of a fluorogram of [<sup>35</sup>S]methionine labeled cellular proteins from human AMA cells (master database) (20). Images can be displayed either in black and white (resembling the original fluorograms) or in color (other images in Fig. 2), depending on the need. As shown in Fig. 2B, each polypeptide is assigned a number by the computer, which facilitates the entry and retrieval of qualitative and quantitative information for any given spot in the gel (20). The standard image can be matched automatically by the computer to other standard or reference gels (Fig. 2C; matching of AMA cellular proteins [left] to MRC-5 proteins [right]) provided a few landmark spots are given manually as reference (indicated with a + in Fig. 2C) to initiate the process.

<sup>2</sup>Abbreviations: CCD, charge couple device; PCNA, proliferating cell nuclear antigen; HPLC, high performance liquid chromatography.



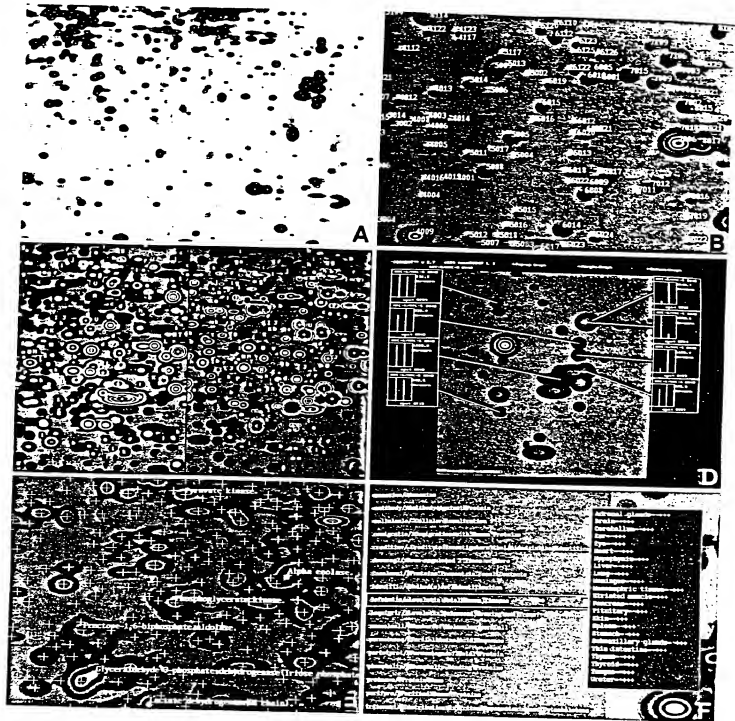
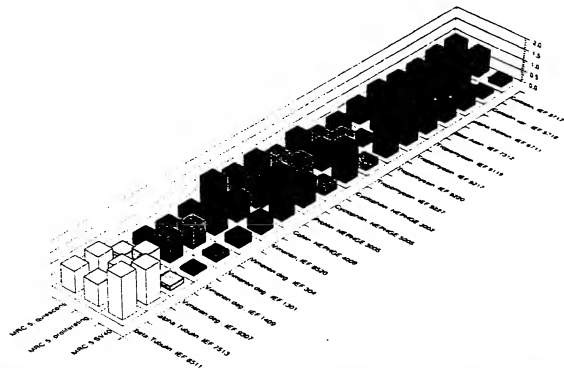
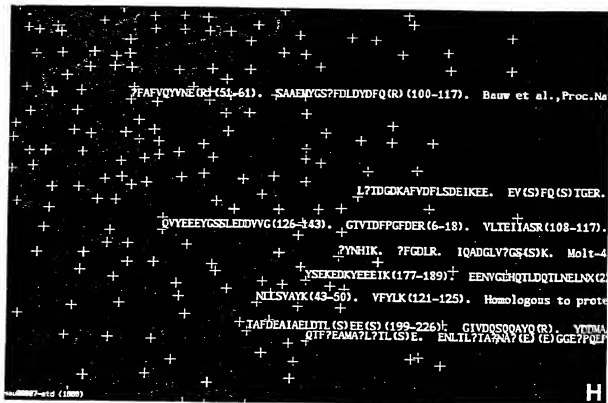


Figure 2. A) Synthetic image of a fraction of an IEF gel of the master image of AMA cellular proteins. B) As in A but showing numbers assigned to each spot. C) Comparison of AMA (left) and normal human embryonal lung MRC-5 fibroblasts (right) IEF protein patterns. Matched proteins are indicated by a + or by the same letters in both gels. Once a protein is matched, information contained in the various categories available in the master AMA database can be transferred. D) Synthetic image of a fraction of an IEF fluorogram of  $^{125}$ I-methionine and SV40 transformed MRC-5 fibroblasts. The histograms show levels of synthesis of a few proteins in MRC-5 (left) and SV40 transformed MRC-5 (right bar) fibroblasts. E) Polypeptides that contain information under the category glycolytic pathway. F) The function peruse annotation for spot allows the operator to inquire about categories and information available for a given protein. G) Relative abundance of cytoskeletal and cytoskeletal-related proteins in quiescent, proliferating, and SV40-transformed MRC-5 fibroblasts. H) Polypeptides that contain information under the category partial amino acid sequences.



**G**



H

The automatic matching process that has been described in detail by Garrels et al. (12) takes about 5 min. Matched proteins are indicated with the same letters in both gels (Fig. 2C). The usefulness of this function is emphasized by the fact that data accumulated on common household proteins can be easily transferred to any other human cellular cell type whose 2-dimensional gel cellular protein pattern is matched

to our standard AMA 2-dimensional gel protein image. Alternatively, if the standard gel is part of a matchset (set of gels in a given experiment) it can be used as a linker gel to compare, for example, the quantitative values of a given protein throughout the experiment (see Fig. 2D; levels of some proteins in normal and SV40 transformed human MRC-5 fibroblasts) or with other standard images in different sets of

cross-matched experiments (18, 22).

Once a standard map of a given protein sample is made, one can enter qualitative annotations to make a reference database. Our master 2-dimensional gel database of transformed human amnion cell (AMA) proteins (20) lists 3430 polypeptides of which 2592 correspond to cellular components, having pI's ranging from 4 to 13 and molecular weights between 8.5 and 230 kDa. The most abundant proteins in the database correspond to total actin (3.87% of total protein; about 90 million molecules per cell) while the lesser abundant of the recorded polypeptides are present in the vicinity of 5000 molecules per cell. Some annotation categories we are using to establish the master AMA database include: 1) protein identification (comigration with purified proteins, 2-dimensional immunoblotting, microsequencing); 2) amounts (total amounts and levels of synthesis); 3) subcellular localization (nuclear, cytoskeletal, membrane, membrane receptors, specific organelles, etc.); 4) antibodies; 5) posttranslational modifications (phosphorylation, glycosylation, methylation etc.); 6) microsequencing; 7) cell cycle specificity (specific variations in levels of synthesis and amount); 8) regulatory behavior (effect of hormones, growth factors, heat shock, etc.); 9) rate of synthesis in normal and transformed cells (proliferation sensitive proteins, cell cycle specific proteins, oncogenes, components of the pathway (or pathways) that control cell proliferation); 10) function (mainly from comigration with proteins of known function); 11) sets of proteins that are coordinately regulated (hierarchy of controls, differential gene expression in various cells, etc.); 12) cDNAs (cloned cDNAs); 13) proteins that are specific to a given disease (systematic comparison of protein patterns of fibroblast proteins from healthy and diseased individuals); 14) expression and exploitation of transfected cDNAs; 15) pathways (metabolic, others); 16) gene localization (genetic and physical); 17) effect of microinjected antibody on patterns of protein synthesis; and 18) secreted proteins.

Information entered for any spot in a given annotation category can be easily retrieved by asking the computer to display the information on the color screen. For example, Fig. 2E shows a synthetic image of a NEPHGE gel (master AMA database) displaying the information contained under the entry glycolytic pathway. Alternatively, one can use the function peruse annotations for spot to directly ask the computer to list all the entries available for a particular protein. By clicking the mouse in a given entry (in this case, presence in fetal human tissues) it is possible to take a quick look at the information in that particular entry (Fig. 2F).

A major obstacle encountered in building comprehensive 2-dimensional gel protein databases is identifying the large number of proteins separated by this technology. In our databases (20, 21), known proteins are identified by one or a combination of the following procedures: 1) comigration with known proteins, 2) 2-dimensional gel immunoblotting using specific antibodies, and 3) microsequencing of Coomassie Brilliant Blue stained human proteins recovered from dried 2-dimensional gels (see next section). Protein identification by means of microsequencing may be difficult, as individual protein members of families with short peptide differences may escape detection. In the gene-protein database of *E. coli* K-12 (14, 23), another major 2-dimensional gel database available at present, proteins are being identified by a wider range of tests that include comigration with purified proteins; genetic criterion (deletion, insertion, frameshift, nonsense, missense, regulatory), plasmid-bearing strains and in vitro synthesis of protein; selective labeling (methylation, phosphorylation); peptide map similarity; and physiological criterion and selective derivatization.

So far we have received nearly 550 antibodies from laboratories all over the world and these are being systematically tested by 2-dimensional gel immunoblotting for antigen determination. Similarly, purified proteins and organelles provided by several laboratories have greatly aided identification of unknown proteins (20, 21). We routinely request antibodies and protein samples and promise the donors to make available all the information we may have accumulated on that particular protein. For example, Table 1 lists entries available for Lipocortin V (IEF SSP 8216), also known as annexin V, VAC- $\alpha$ , endonexin II, renocortin, chromobindin-5', anticoagulant protein, PAP-I,  $\gamma$ -calicmedin, IBC, calphobindin, and anchorin CII.

As mentioned previously, one distinct advantage of 2-dimensional gel electrophoresis is the possibility of studying quantitative variations in cellular protein patterns that may lead to identification of groups of proteins that are expressed coordinately during a given biological process. Quantitation, however, is not an easy task as reflected by the lack of published data on global cellular protein patterns. We believe this is partly due to difficulties in obtaining sets of gels that are suitable for computer analysis (streaking, material remaining at the origin, etc.) as well as to limitations (laborious editing time, need of calibration strips to merge images, limited dynamic range, etc.) in the computer analysis systems available at the moment. Perhaps the most advanced quantitative studies published so far using computer analysis have been carried out by Garrels and co-workers (18, 22). In particular, these investigators have established a quantitative rat protein database (18, 22) designed to study growth control (proliferation, growth inhibitors, and stimulation) and transformation in well-defined groups of cell lines obtained by transformation of rat REF52 cells with SV40, adenovirus, and the Kirsten murine sarcoma virus. These studies have revealed clusters of proteins induced or repressed during growth to confluence as well as groups of transformation-sensitive proteins that respond in a differential fashion to transformation by DNA and RNA viruses. A most interesting feature of this quantitative database is the discovery of a group of coregulated proteins that show similar expression patterns as the cell cycle-regulated DNA replication protein known as proliferating cell nuclear antigen (PCNA)/cyclin (45).

In our human databases, most quantitations have been carried out by estimating the radioactivity contained in the polypeptides by direct counting of the gel pieces in a scintillation counter (20, 21). Up to 700 proteins can be cut out through appropriate exposed films in a period of time comparable to that required for editing a synthetic image. Manual quantitation of this large number of spots is difficult without the assistance of a master reference image and a numbering system that can be used to identify the spots. Using this approach, we have recorded quantitative changes in the relative abundance of 592 [ $^{35}$ S]methionine-labeled proteins synthesized by quiescent, proliferating, and SV40 transformed human embryonic lung MRC-5 fibroblasts (21). Some data concerning cytoskeletal and cytoskeletal-related proteins are presented in Fig. 2G. Our studies as well as those of Garrels and co-workers (18, 22) may in the long run help define patterns of gene expression that are characteristic of the transformed state.

## OTHER 2-DIMENSIONAL GEL PROTEIN DATABASES

As mentioned previously there are other 2-dimensional gel databases available in computer form that have been pub-

TABLE 1. Some entries for lipocortin V in the human A.M.A. 2-dimensional gel protein database

Entries for lipocortin V (IEF SSP 8216)	Information entered
1. Protein name	Lipocortin V, renocortin, chromobindin-5, endonexin I, anticoagulant protein, PAP-I, VAC- $\alpha$ , 35- $\gamma$ -calcimedlin, IBC, calphobindin I, anchorn CII, annex V
2. Percentage of total protein	0.1105 (about 2,800,000 molecules per cell)
3. Apparent molecular weight (mr)	33.3 kDa
4. Isoelectric point (pl)	4.76
5. Method (or methods) of identification	Microsequencing, 2-dimensional immunoblotting, Comigration
6. Credit to investigators that aided in identification	G. Bauw, J. Vandekerckhove, and colleagues, Rijksuniversiteit Gent; B. Pepinsky, BIOGEN, Cambridge; N.G. Ahn, University of Washington
7. Antibody against protein	Polyclonal (rabbit, antibody no. 20), B. Pepinsky, BIOGEN, Cambridge
8. Comigration with human proteins	Lipocortin V, N.G. Ahn, Howard Hughes Medical Institute, Washington University
9. Cellular localization	Subcortical membrane
10. Calcium-phospholipid-dependent membrane proteins	Lipocortin V
11. Function	Regulation of various aspects of inflammation, immune response, blood coagulation and differentiation
12. Partial amino acid sequence	GTVTDGPGFDER (7-18), VLTEIHASR (109-117), QVYEEYGVSSLEDIVVG (127-143), ?GTDEENFITIFGT(R) (187-201)
13. cDNA sequence	Known: R. Blake et al., <i>J. Biol. Chem.</i> 263, 10799-10811, 1988 (pl = 4.76 from translated sequence)
14. Levels in fetal human tissues	Adrenal glands = - - - - ; brain = - - - - ; cerebellum = - - - - ; ear = - - - - ; eye = - - - - ; heart = - - - - ; hypophysis = - - - - ; liver = - - - - ; lung = - - - - ; meninges = - - - - ; mesonephric tissue = - - - - ; striated muscle = - - - - ; pancreas = - - - - ; skin = - - - - ; spleen = - - - - ; stomach = - - - - ; submandibular gland = - - - - ; small intestine = - - - - ; thymus = - - - - ; thyroid gland = - - - - ; tongue = - - - - ; ureter = - - - -
15. Levels in quiescent, proliferating, and transformed MRC-5 fibroblasts	Q (quiescent) = 1.1; P (proliferating) = 1.0; T (SV40 transformed) = 0.3
16. Distribution in Triton supernatant and cytoskeletons	Mainly supernatant

lished in extenso: these correspond to the *E. coli* K-12 protein-gene database (14, 23) and to the rat REF52 database (18, 22).

The *E. coli* K-12 cellular protein-gene database is perhaps the most complete of all databases reported so far and eventually it should trace each protein back to its structural gene. Information contained in this database includes: gene/protein name (protein name, EC number, gene name); 2-dimensional gel spot designations (x-y coordinates from reference gels, alphanumeric designation); genetic information (linkage map location, physical map location, Genebank code, sequence reference, location on Kohara clones); biochemical information (molecular weight, pl, number of residues of each amino acid, mole percent of each amino acid, total number of amino acids in a polypeptide), and regulatory information (cellular level of protein in different media and different temperature, member of regulon, member of stimulon). Major advances of this database are envisaged in the future in view of the eminent sequencing of

the whole *E. coli* genome as well as the development of improved methods to express cloned genes.

The rat REF52 2-dimensional gel protein database lists about 1600 proteins that have been recorded using the QUEST analysis system (18, 22). Included in this quantitative database are 1) protein names (cytoskeletal and heat shock proteins as well as various nuclear, mitochondrial, and cytoplasmic proteins), 2) annotations (subcellular localization, modification, recognition by specific antibodies, coprecipitation, NH<sub>2</sub>-terminal sequence, cross-reference to protein sequence information and references to the literature), 3) protein sets (cytoskeletal proteins, phosphoproteins, sets of proteins with PCNA/cyclin-like properties, etc.) and 4) general quantitative data (protein synthesis during growth of normal REF52 cells to confluence and quiescence, and after restimulation of growth-inhibited cells).

In addition to the 2-dimensional gel databases mentioned so far there are several smaller cellular databases being established in human (normal human diploid fibroblasts, lymphocytes, etc.).

phocytes, leukocytes, leukemic cells) mouse (NIH/3T3 cells, T lymphocytes), *Apisya*, yeast (*Saccharomyces cerevisiae*), plants (wheat, barley, sorghum), and *Euglena*. Databases of tissue protein, (brain, whole mouse, liver) and body fluid proteins (plasma proteins, cerebrospinal fluid, urine, and milk) are being established in several laboratories. The reader is directed to the review by Celis et al. (4) for details and references concerning these databases.

#### MICROSEQUENCING HAS ADDED A NEW DIMENSION TO COMPREHENSIVE 2-DIMENSIONAL GEL DATABASES: A DIRECT LINK BETWEEN PROTEINS AND GENES

The development of highly sensitive amino acid gas-phase or liquid-phase sequencers (24), together with the establishment of efficient protein and peptide sample preparation methods, has opened the possibility to perform a systematic sequence analysis of proteins resolved by 2-dimensional gel electrophoresis. Indeed, generated pieces of protein sequences can be used to search for protein identity (comparison with available sequences stored in databanks) as well as for preparing specific DNA probes for cloning of as yet uncharacterized proteins (Fig. 1). In addition, partial protein sequences can be stored in 2-dimensional gel databases (for example, see Fig. 2H) and offer a unique link between proteins and genes (Fig. 1).

In the early 1970s gel electrophoresis was used to purify proteins for sequencing purposes (reviewed by Weber and Osborn in ref 25). Proteins were recovered by diffusion and sequenced by the manual dansyl-Edman degradation at the nanomole level. This technique was further refined by using electro-elution to recover proteins and by miniaturizing the system (26). This method has been used extensively, but showed increasing drawbacks (low yields, protein samples contaminated by free amino acids, and  $\text{NH}_2$ -terminal blocking) as the amounts of handled protein gradually became smaller (e.g., at the 10 picomol level).

Most of the problems referred to above have been minimized with the introduction of protein-electroblotting procedures (27-32). When proteins are blotted on chemically inert membranes, it is possible to sequence the immobilized proteins directly without additional manipulations. Thus, depending on the amount of bound protein and its nature, this direct sequencing procedure generally yields  $\text{NH}_2$ -terminal sequences containing 10-40 residues. As such, this technique was used to identify, by their  $\text{NH}_2$ -terminal sequences, differentially expressed major proteins from total cellular extracts separated on 2-dimensional gels. A major difficulty encountered in this procedure is the occurrence of frequent artefactual blockage of the proteins. Several studies suggest that this phenomenon is mainly due to reaction with contaminants (particularly unpolymerized acrylamide present in the gel) and to a high dilution of the protein (low concentration of the protein per unit membrane surface). In addition to this primarily technical problem, many proteins are blocked in vivo by acylation or by a pyrrolidone carboxylic acid cap.

The problem of partial or complete  $\text{NH}_2$ -terminal blockage can be circumvented by generating internal amino acid sequences. This is achieved by fragmenting the protein present in the gel (gel in situ cleavage) or by cleaving it while bound to the membrane (membrane in situ cleavage) (33-35). In both cases, proteins are either cleaved in a restricted way (e.g., by limited enzymatic digestion or by using restriction chemical cleavage conditions) or fragmented into smaller peptides.

Of the different combinations examined, we had good results by using exhaustive proteolytic digestion on membrane-immobilized proteins. This method has been described for Ponceau red-stained proteins on nitrocellulose blots (34), for Amido-black-stained Immobilized-bound proteins, and for fluorescamine-detected proteins on glass fiber membranes (35). The proteases used (trypsin, chymotrypsin, or pepsin) cleave at multiple sites, generating small peptides that elute from the blot into the digestion buffer from which they are purified by reversed-phase high performance liquid chromatography (HPLC) before being sequenced individually. Although each of these manipulations could be expected to result in a reduced yield of final sequence information, we were surprised that the peptides could be sequenced with high efficiency. In our hands, this approach could be routinely applied to gel-purified proteins available in amounts ranging from 5 to 10  $\mu\text{g}$ , and often yielded sequence information covering more than 30% of the total protein. As membrane-immobilized proteins are not homogeneously digested, but rather show protease sensitivity next to resistant regions, the number of peptides generated is much lower than expected from the number of potential cleavage sites. Consequently, HPLC peptide chromatograms are less complex and most peptides can be recovered in pure form.

As only limited amounts of a protein mixture can be loaded on a 2-dimensional gel, proteins of interest are often obtained in yields insufficient for the currently available sequencing technology. More material can be obtained by enriching for a certain subcellular fraction (purified cell organelles) or by exploiting affinity (dyes, metals, drugs, etc) or hydrophobic properties of proteins before gel analysis. All of the sequencing results accumulated so far in the human protein database (20) (a few are shown in Fig. 2H) have been obtained from analysis of protein spots collected from 2-dimensional gels that had been stained with Coomassie blue according to standard procedures and dried for storage. Proteins are recovered from the collected gel pieces by a protein-elution-concentration device, combined with gel electrophoresis and electroblotting. Details of this technique have been reported in a previous communication (42) and a brief outline is given below.

Combined gel pieces are allowed to swell in gel sample buffer (a total volume of 1.5 ml). The gel pieces combined with the supernatant are then collected into a large slot made in a new gel. The slot is further filled with Sephadex G-10 equilibrated in gel sample buffer. During consecutive gel electrophoresis, most of the electrical current passes on the side of the slot instead of passing through the slot. This results in both a vertical stracking and horizontal contraction of the protein band. With this device the protein is efficiently eluted from the gel pieces and concentrated from a large volume into a narrow spot. The highly concentrated (about 5  $\text{mm}^2$ ) protein spot is then electroblotted on PVDF-membranes, stained with Amide black, and in situ digested with trypsin. The peptides generated during digestion elute from the membrane into the supernatant, and can be separated by narrow bore reversed-phase HPLC and collected individually for sequence analysis.

Using this and previous procedures (37, 39, 42), we have so far analyzed 70 protein spots collected from 2-dimensional gels (20, and unpublished observations) (see for example Fig. 2H). The sequence information amounts to 2100 allocated residues corresponding to an average of 30 residues per protein spot. So far we have made cDNAs of many of the unknown proteins that have been microsequenced, and a substantial number has been cloned and sequenced. All available information indicates that it may be possible to obtain partial sequence information from most of

the proteins that can be visualized by Coomassie Brilliant Blue staining.

Partial protein sequences are stored in the database as displayed in Fig. 2H, and it should be possible in the near future to interface this information with forthcoming DNA sequence data from the human genome project. In the long run, as the human genome sequences become available it will be possible to assign partial protein sequences to genes for which the full DNA sequence and chromosomal location are known (Fig. 1).

## SUMMARY

The studies presented in this brief review are intended to demonstrate the usefulness of computer-aided 2-dimensional gel electrophoresis and microsequencing to analyze cellular protein patterns, and to link protein and DNA information. As more information is gathered worldwide, comprehensive databases will depict an integrated picture of the expression levels and properties of the thousands of proteins that orchestrate most cellular functions.

Clearly, databases allow easy access to a large body of data and provide an efficient medium to communicate standardized protein information. In the future, databases will foster a wide variety of biological information that can be used to support collaborative research projects in basic and applied biology as well as in clinical research (2, 5, 46). Once a protein is identified in a particular database all the information gathered on it can be made available to the scientist. However, many problems must be solved before protein databases become of general use to the scientific community. A most urgent one is to promote standardization of the gel running conditions so that data produced in a given laboratory may be used worldwide. Surprisingly, the gel running technology as it stands today is still a craftsmanship art.

Finally, comprehensive, computerized databases of proteins, together with recently developed techniques to microsequence proteins, offer a new dimension to the study of genome organization and function (Fig. 1). In particular, human protein databases may become increasingly important in view of the concerted effort to map and sequence the entire human genome. This formidable task is expected to dominate biological research in the next decades. [F]

We would like to thank S. Himmelstrup-Jørgensen for typing the manuscript and O. Sønderskov for photography. Work in the authors' laboratories was supported by grants from the Danish Biotechnology Programme, the Danish Cancer Foundation, and the Commission of the European Communities.

## REFERENCES

- O'Farrell, P. H. (1975) High resolution two-dimensional electrophoresis of proteins. *J. Biol. Chem.* **250**, 4007-4021.
- Special Issue: Two-dimensional gel electrophoresis. *Clin. Chem.* **28**, 1982.
- Celis, J. E., and Bravo, R., eds. (1984) *Two-Dimensional Gel Electrophoresis of Proteins: Methods and Applications*. Academic, New York.
- Celis, J. E., Madsen, P., Gesser, B., Kwee, S., Nielsen, H. V., Rasmussen, H. H., Honoré, B., Leffers, H., Ratz, G. P., Basse, B., Lauridsen, J. B., and Celis, A. (1989) Protein databases derived from the analysis of two-dimensional gels. In *Advances in Electrophoresis* (Chrambach, C., ed) VCH, Weinheim, Germany.
- Special Issue: Two-dimensional gel electrophoresis in cell biology. (Celis, J. E., ed) *Electrophoresis* **11**, 1990.
- Celis, J. E., Honoré, B., Bauw, G., and Vandekerckhove, J. (1990) Comprehensive computerized 2D gel protein databases offer a global approach to the study of the mammalian cell. *BioEssays* **12**, 93-98.
- Garrels, J. I. (1983) Two-dimensional gel electrophoresis and computer analysis of proteins synthesized by cloned cell lines. *Mol. Cell. Biochem.* **100**, 411-423.
- Anderson, N. L., Hoimann, J. P., Gemmel, A., and Taylor, S. (1984) Global approaches to the quantitative analysis of gene expression patterns observed by two-dimensional gel electrophoresis. *Clin. Chem.* **30**, 2031-2036.
- Garrels, J. I., Farrar, J. T., and Burwell, C. B. (1984) The Quest system for computer-analyzed two-dimensional electrophoresis of proteins in *Two-Dimensional Gel Electrophoresis of Proteins: Methods and Applications* (Celis, J. E., and Bravo, R., eds) pp 37-91. Academic, New York.
- Vincens, P., and Tarroux, P. (1988) Two-dimensional electrophoresis computerized processing. *Int. J. Biochem.* **20**, 499-509.
- Appel, R., Hochstrasser, D., Roch, C., Funk, M., Müller, A. F., and Pellegrini, C. (1988) Automatic classification of two-dimensional gel electrophoresis pictures by heuristic clustering analysis: a step toward machine learning. *Electrophoresis* **9**, 136-142.
- Lemkin, P. F., and Lester, E. P. (1989) Database and search techniques for two-dimensional gel protein data: a comparison of paradigms for exploratory data analysis and prospects for biological modeling. *Electrophoresis* **10**, 122-140.
- Miller, M. J. (1989) Computer-assisted analysis of two-dimensional gel electrophoretograms. *Adv. Electrophoresis* **3**, 182-217.
- Phillips, T. D., Vaughn, V., Bloch, P. L., and Neidhardt, F. C. (1987) In *Escherichia coli and Salmonella typhimurium. Cellular and Molecular Biology. Gene-Protein Index of Escherichia coli K-12*, 2. ed. (Neidhardt, F. C., Ingraham, J. I., Low, K. B., Magasanik, B., Schaechter, M., and Umberger, H. E., ed) pp 919-966. American Society for Microbiology, Washington, D.C.
- Celis, J. E., Ratz, G. P., Celis, A., Madsen, P., Gesser, B., Kwee, S., Madsen, P. S., Nielsen, H. V., Yde, H., Lauridsen, J. B., and Basse, B. (1988) Towards establishing comprehensive databases of cellular proteins from transformed human epithelial amnion cells (AMA) and normal peripheral blood mononuclear cells. *Leukemia* **2**, 561-601.
- Special Issue: Protein databases in two-dimensional electrophoresis. (Celis, J. E., ed) *Electrophoresis* **2**, 1989.
- Celis, J. E., Ratz, G. P., Madsen, P., Gesser, B., Lauridsen, J. B., Brogaard-Hansen, K. P., Kwee, S., Rasmussen, H. H., Nielsen, H. V., Crüger, D., Basse, B., Leffers, H., Honoré, B., Möller, O., and Celis, A. (1989) Computerized, comprehensive databases of cellular and secreted proteins from normal human embryonic lung MRC-5 fibroblasts: identification of transformation and/or proliferation sensitive proteins. *Electrophoresis* **10**, 76-115.
- Garrels, J. I., and Franza, B. R. (1989) The REF52 protein database. Methods of database construction and analysis using the Quest system and characterizations of protein patterns from proliferating and quiescent REF52 cells. *J. Biol. Chem.* **264**, 5283-5298.
- Celis, J. E., Crüger, D., Küll, J., Dejgaard, K., Lauridsen, J. B., Ratz, G. P., Basse, B., Celis, A., Rasmussen, H. H., Bauw, G., and Vandekerckhove, J. (1990) A two-dimensional gel protein database of noncultured total normal human epidermal keratinocytes: identification of proteins strongly up-regulated in psoriatic epidermis. *Electrophoresis* **11**, 242-254.
- Celis, J. E., Gesser, B., Rasmussen, H. H., Madsen, P., Leffers, H., Dejgaard, K., Honoré, B., Olsen, E., Ratz, G., Lauridsen, J. B., Basse, B., Mouritzen, S., Hellerup, M., Andersen, A., Walbum, E., Celis, A., Bauw, G., Puyve, M., Van Damme, J., and Vandekerckhove, J. (1990) Comprehensive two-dimensional gel protein databases offer a global approach to the analysis of human cells: the transformed amnion cells (AMA) master database and its link to genome DNA sequence data. *Electrophoresis* **12**, 989-1071.

21. Celis, J. E., Dejgaard, K., Madsen, P., Leffers, H., Gesser, B., Honoré, B., Rasmussen, H. H., Olsen, E., Lauridsen, J. B., Ratz, G., Mouritzen, S., Hellerup, M., Andersen, A., Walbum, E., Celis, A., Bauw, G., Puype, M., Van Damme, J., and Vandekerckhove, J. (1990) The MRC-5 human embryonal lung fibroblast two-dimensional gel cellular protein database: quantitative identification of polypeptides whose relative abundance differs between quiescent, proliferating and SV40 transformed cells. *Electrophoresis* 12, 1072-1113.
22. Garrels, J. I., Franza, B. R., Chang, C., and Luster, G. (1990) Quantitative exploration of the REF52 protein database: cluster analysis reveals the major protein expression profiles in responses to growth regulation, serum stimulation, and viral transformation. *Electrophoresis* 12, 1114-1130.
23. Van Bozelen, R. A., Hutton, M. E., and Neidhardt, F. C. (1990) Gene-protein database of *Escherichia coli* K-12, 3rd ed. *Electrophoresis* 12, 1131-1166.
24. Hewick, R. M., Hunkapiller, M. W., Hood, L. E., and Drever, W. J. (1981) A gas-liquid solid phase peptide and protein sequencer. *J. Biol. Chem.* 256, 7990-7997.
25. Weber, K., and Osborn, M. (1985) In *The Proteins and Sodium Dodecyl Sulfate, Molecular Weight Determination on Polyacrylamide Gels and Related Procedures* (Neurath, H. et al., eds) Vol. 1, pp. 179-223. Academic, New York.
26. Hunkapiller, M. W., Lujan, E., Ostrander, F., and Hood, L. E. (1983) Isolation of microgram quantities of proteins from polyacrylamide gels for amino acid sequence analysis. *Methods Enzymol.* 91, 227-236.
27. Vandekerckhove, J., Bauw, G., Puype, M., Van Damme, J., and Van Montagu, M. (1985) Protein-blotting on polybrene-coated glass-fiber sheets. *Eur. J. Biochem.* 152, 9-19.
28. Aebersold, R. H., Teplow, D. B., Hood, L. E., and Kent, S. B. H. (1986) Electrophoretic onto activated glass. *J. Biol. Chem.* 261, 4229-4238.
29. Bauw, G., De Loose, M., Inzé, D., Van Montagu, M., and Vandekerckhove, J. (1987) Alterations in the phenotype of plant cells studied by NH<sub>2</sub>-terminal amino acid-sequence analysis of proteins electrophoretically separated from two-dimensional gel-separated total extracts. *Proc. Natl. Acad. Sci. USA* 84, 4806-4810.
30. Matsuda, P. (1987) Sequence from picomole quantities of proteins electrophoretically separated from polyvinylidene difluoride membranes. *J. Biol. Chem.* 262, 10033-10038.
31. Eckerskorn, C., Mewes, W., Goretzki, H., and Lottspeich, F. (1985) A new siliconized-glass fiber as support for electrophoretic analysis of electrophoretically separated proteins. *Eur. J. Biochem.* 176, 509-519.
32. Moos, M. Jr., Nguyen, N. Y., and Liu, T.-Y. (1988) Reproducible high yield sequencing of proteins electrophoretically separated and transferred to an inert support. *J. Biol. Chem.* 263, 6003-6008.
33. Kennedy, T. E., Gawinowicz, M. A., Barzilai, A., Kandel, E. R., and Sweet, J. D. (1988) Sequencing of proteins from two-dimensional gels by using in situ digestion and transfer of peptides to polyvinylidene difluoride membranes: application to protein associated with sensitization in *Apis mellifera*. *Proc. Natl. Acad. Sci. USA* 85, 7008-7012.
34. Aebersold, R. H., Leavitt, J., Saavedra, R. A., Hood, L. E., and Kent, S. B. H. (1987) Internal amino acid sequence analysis of protein separated by one- or two-dimensional gel electrophoresis after in situ protease digestion on nitrocellulose. *Proc. Natl. Acad. Sci. USA* 84, 6870-6872.
35. Bauw, G., Van Den Bulcke, M., Van Damme, J., Puype, M., Van Montagu, M., and Vandekerckhove, J. (1988) Protein electroblotting on polybrene-coated glassfiber and polyvinylidene difluoride membranes: an evaluation. *J. Prot. Chem.* 7, 194-198.
36. Celis, J. E., Ratz, G., Madsen, P., Gesser, B., Lauridsen, J. B., Leffers, H., Rasmussen, H. H., Nielsen, H. V., Cruet, D., Basse, B., Honoré, B., Möller, O., Celis, A., Vandekerckhove, J., Bauw, G., Van Damme, J., Puype, M., and Van Den Bulcke, M. (1989) Comprehensive, human cellular protein databases and their implication for the study of genome organization and function. *FEBS Lett.* 244, 247-254.
37. Bauw, G., Van Damme, J., Puype, M., Vandekerckhove, J., Gesser, B., Lauridsen, J. B., Ratz, G. P., and Celis, J. E. (1989) Protein-electroblotting and -microsequencing strategies in generating protein databases from two-dimensional gels. *Proc. Natl. Acad. Sci. USA* 86, 7701-7705.
38. Aebersold, R., and Leavitt, J. (1990) Sequence analysis of proteins separated by polyacrylamide gel electrophoresis. Towards an integrated protein database. *Electrophoresis* 11, 517-527.
39. Bauw, G., Rasmussen, H. H., Van Den Bulcke, M., Van Damme, J., Puype, M., Gesser, B., Celis, J. E., and Vandekerckhove, J. (1990) Two-dimensional gel electrophoresis, protein electroblotting and microsequencing: a direct link between proteins and genes. *Electrophoresis* 11, 528-536.
40. Tempst, P., Link, A. J., Riviere, L. R., Fleming, M., and Ellicott, C. (1990) Internal sequence analysis of protein separated on polyacrylamide gels at the submicrogram level: improved methods, applications and gene cloning strategies. *Electrophoresis* 11, 537-553.
41. Eckerskorn, C., and Lottspeich, F. (1990) Combination of two-dimensional gel electrophoresis with microsequencing and amino acid composition analysis: improvement of speed and sensitivity in protein characterization. *Electrophoresis* 11, 554-561.
42. Rasmussen, H. H., Van Damme, J., Bauw, G., Puype, M., Gesser, B., Celis, J. E., and Vandekerckhove, J. (1991) In *Methods in Protein Sequence Analysis* (Jornvall, H., and Höög, J. O., eds) pp. 103-114. Eighth International Conference on Methods in Protein Sequence Analysis. Birkhäuser Verlag, Boston.
43. Olson, A. D., and Miller, M. J. (1988) Elsevier: quantitative computer analysis of sets of two-dimensional gel electrophoretograms. *Anal. Biochem.* 169, 49-70.
44. Vincens, P., Paris, N., Pujol, J. L., Gaboriaud, C., Rabilloud, T., Pennecler, J., Natherat, P., and Tarroux, P. (1986) HERMES: a second generation approach to the automatic analysis of two-dimensional electrophoresis gels. Part I: Data acquisition. *Electrophoresis* 7, 347-356.
45. Celis, J. E., Madsen, P., Celis, A., Nielsen, H. V., and Gesser, B. (1987) Cyclin (PCNA, auxiliary protein of DNA polymerase  $\delta$ ) is a central component of the pathway(s) leading to DNA replication and cell division. *FEBS Lett.* 220, 1-7.
46. Anderson, N. G., and Anderson, N. L. (1982) The human protein index. *Clin. Chem.* 28, 739-748.

the 1990s, the number of people in the world who are under 15 years of age is expected to increase from 1.1 billion to 1.5 billion.

As the world's population grows, the demand for food and other resources will increase. This will put pressure on the environment and on the world's food supply. It is important that we find ways to meet this demand without harming the environment.

One way to do this is to use sustainable agriculture. This means using farming methods that do not harm the environment and that can be used over and over again. Sustainable agriculture can help us to meet the world's growing demand for food without harming the planet.

Another way to do this is to use renewable resources. These are resources that can be replaced naturally, such as wind, water, and solar energy. Using renewable resources can help us to meet the world's growing demand for energy without depleting the planet's resources.

Finally, we can help to meet the world's growing demand for food and other resources by using less of them. This means using things more carefully and recycling whenever possible. By using less, we can help to reduce the pressure on the environment and on the world's food supply.

There are many other ways that we can help to meet the world's growing demand for food and other resources. The important thing is that we all do our part. By working together, we can make sure that the world has enough food and other resources for everyone.

One of the most important things we can do is to educate people about the need to use resources sustainably. This means teaching people about the importance of the environment and about the need to use resources carefully. By educating people, we can help to create a world where everyone is responsible for the planet.

Another important thing we can do is to support sustainable agriculture. This means buying products from farmers who use sustainable farming methods. By supporting sustainable agriculture, we can help to ensure that the world's food supply is sustainable.

Finally, we can help to meet the world's growing demand for food and other resources by using renewable resources. This means using energy from wind, water, and solar power. By using renewable resources, we can help to ensure that the world's energy supply is sustainable.

There are many other things that we can do to help meet the world's growing demand for food and other resources. The important thing is that we all do our part. By working together, we can make sure that the world has enough food and other resources for everyone.

One of the most important things we can do is to use less of the resources that we need. This means using things more carefully and recycling whenever possible. By using less, we can help to reduce the pressure on the environment and on the world's food supply.

Another important thing we can do is to support sustainable agriculture. This means buying products from farmers who use sustainable farming methods. By supporting sustainable agriculture, we can help to ensure that the world's food supply is sustainable.

Finally, we can help to meet the world's growing demand for food and other resources by using renewable resources. This means using energy from wind, water, and solar power. By using renewable resources, we can help to ensure that the world's energy supply is sustainable.



Bengt Bjellqvist\*  
 Bodil Basse  
 Eydfinnur Olsen  
 Julio E. Celis

Institute of Medical Biochemistry  
 and Danish Centre for Human  
 Genome Research, Aarhus  
 University, Aarhus

## Reference points for comparisons of two-dimensional maps of proteins from different human cell types defined in a pH scale where isoelectric points correlate with polypeptide compositions

A highly reproducible, commercial and nonlinear, wide-range immobilized pH gradient (IPG) was used to generate two-dimensional (2-D) gel maps of [<sup>35</sup>S]methionine-labeled proteins from noncultured, unfractionated normal human epidermal keratinocytes. Forty-one proteins, common to most human cell types and recorded in the human keratinocyte 2-D gel protein database were identified in the 2-D gel maps and their isoelectric points (pI) were determined using narrow-range IPGs. The latter established a pH scale that allowed comparisons between 2-D gel maps generated either with other IPGs in the first dimension or with different human protein samples. Of the 41 proteins identified, a subset of 18 was defined as suitable to evaluate the correlation between calculated and experimental pI values for polypeptides with known composition. The variance calculated for the discrepancies between calculated and experimental pI values for these proteins was 0.001 pH units. Comparison of the values by the *t*-test for dependent samples (paired test) gave a *p*-level of 0.49, indicating that there is no significant difference between the calculated and experimental pI values. The precision of the calculated values depended on the buffer capacity of the proteins, and on average, it improved with increased buffer capacity. As shown here, the widely available information on protein sequences cannot, *a priori*, be assumed to be sufficient for calculating pI values because post-translational modifications, in particular N-terminal blockage, pose a major problem. Of the 36 proteins analyzed in this study, 18–20 were found to be N-terminally blocked and of these only 6 were indicated as such in databases. The probability of N-terminal blockage depended on the nature of the N-terminal group. Twenty six of the proteins had either M, S or A as N-terminal amino acids and of these 17–19 were blocked. Only 1 in 10 proteins containing other N-terminal groups were blocked.

### 1 Introduction

As compared with carrier ampholyte isoelectric focusing (CA-IEF), the application of immobilized pH gradients (IPGs) in the first dimension in 2-D gel electrophoresis offers improved reproducibility [1] because the nature of the pH gradient makes the resulting focusing positions insensitive to the focusing time [2] and to the type of sample applied [3]. The recently introduced ready-made IPG strips [4] seem to be an ideal substitute for the carrier ampholyte gradients, which until now have been the most commonly used first dimensions in 2-D gel electrophoresis. The availability of standardized first dimensions opens the possibility of comparing 2-D gel maps of various cell types generated in different laboratories, provided that the focusing positions of a number of easily recognizable polypeptide spots common to the cell types

in question are known. Even though this approach is limited to experiments performed with the same standardized IPG, the flexibility provided by IPGs allows the pH gradient to be adjusted to the requirements of a particular experiment.

Exchange and communication of 2-D gel protein data requires a pH scale that is independent of the particular IPG used and by which the results can be described. The introduction of carbamylation trains and the relation of focusing positions to the spots in these trains represented a step forward towards solving the reproducibility problem experienced with carrier ampholyte focusing [5]. Problems associated with the use of carbamylation trains were mainly due to lack of temperature control and to the use of nonequilibrium focusing conditions. Accordingly, the pattern variation involved not only the resulting pH gradients, but also the relative spot positions as related to each other and to spots in the carbamylation trains. Even though the question of reproducibility has, to a large extent, been solved, the carbamylation trains are still not ideal as markers because the spots in the trains do not represent defined entities but rather a large number of differently carbamylated peptides having close pI values. As a result, the spots are large and poorly defined as compared to the ordinary polypeptide spots in 2-D gel maps.

**Correspondence:** Professor J. E. Celis, Institute of Medical Biochemistry and Danish Centre for Human Genome Research, Aarhus University, DK-8000 Aarhus C, Denmark

**Abbreviations:** CA-IEF, carrier ampholyte-isoelectric focusing; SSP, sample spot number

\* Present address: Pharmacia Biotech AB, S-751 82 Uppsala, Sweden

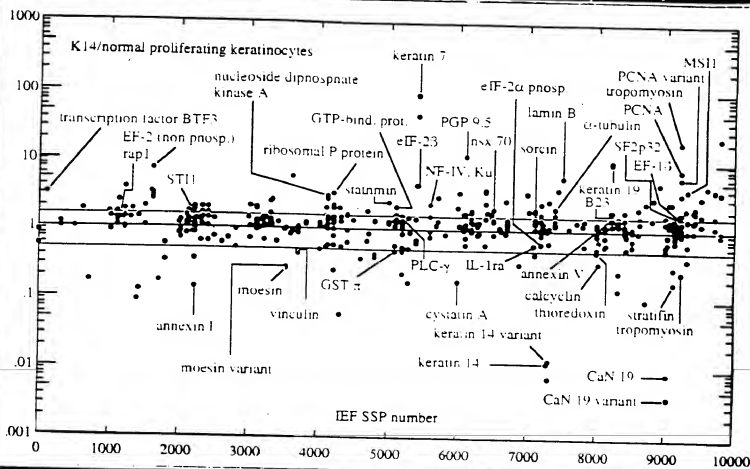
© VCH Verlagsgesellschaft mbH, 0451 Weinheim, 1994

0175-0335/94/0304-0529 \$5.00+25/0

# ELECTROPHORESIS

An International Journal

3-4'94



PAPER SYMPOSIUM

## ELECTROPHORESIS IN CANCER RESEARCH

Guest Editor: Julio E. Celis



ISSN 0173-0835 ELECTDN 15 (3-4) 307-556 (1994) Vol. 15 No. 3-4 March/April 1994

Neidhardt *et al.* [6] defined the pH gradient in 2-D gel experiments by *pI* markers whose *pI* values were calculated from the amino acid composition. Focusing positions of other polypeptides could be predicted from their composition but the *pK* values needed for the *pI* calculations were unknown. Various groups employing this approach do not use the same *pK* values [6, 7] and therefore, the *pI* values derived in this way cannot be expected to describe the variation of the hydrogen ion activity. In spite of this fact, it is still possible to make approximate predictions of focusing positions because the *pK* values used to define the pH gradient are also used to calculate *pI* values and to predict the focusing positions. Errors in *pK* assignments are therefore compensated. A pH scale which correctly reflects the variation in hydrogen ion activity during focusing should improve the precision of the predictions, but this has never been implemented with CA-IEF focusing as a first dimension in 2-D gel electrophoresis. The main reason for this are the problems associated with pH measurements in focused gels containing high concentrations of urea.

IPGs can be described from the concentration variation of the immobilized groups, provided that the *pK* values of these groups are known for the conditions prevailing during focusing. To avoid measurements on gels, Gianazza *et al.* [8] suggested the use of *pK* values derived by addition of determined *pK* shifts. Recently, direct determinations of *pK* differences between immobilized groups in IPGs were made by determining *pI*-*pK* values in overlapping narrow-range IPGs [9, 10] and the results verified the applicability of the Gianazza approach. A description of the focusing results in a pH scale, which correctly describes the variation of the hydrogen ion activity for the focusing conditions used, not only allows the comparison of 2-D gel maps generated with different IPGs, but also opens the possibility of correlating the focusing position of a polypeptide with its composition [9]. Experiments by Bjellqvist *et al.* [9, 10] have implied that pH scales showing good correlation between calculated and experimental *pI* values can be derived for any of the conditions commonly used for focusing in connection with 2-D gel electrophoresis. These pH scales are then defined through the *pK* values of the immobilized groups in the IPG containing gel. To be useful for interlaboratory comparisons, however, the pH scale has to be defined through *pI* values of easily recognizable spots present in the 2-D gel map. So far, *pI* determinations in a useful pH scale, combined with determinations of *pK* values needed for *pI* calculations, have only been made for the pH range 4.5–6.5 at 10°C [9]. CA-IEF focusing as described by O'Farrell [11] does not control the temperature of the first dimension, which can be expected to be slightly above room temperature. With IPGs, the temperature commonly used is about 20°C [4, 12] or 25°C [13] and this is a critical parameter that needs to be controlled [14].

The present work was designed to compare 2-D gel maps of different cell types in a laboratory applying both CA-IEF and IPG focusing at a common temperature. To this end we have generated 2-D gel maps of proteins from noncultured, unfractionated normal human epidermal keratinocytes with IPG in the first dimension

and a focusing temperature of 25°C. We have used commercial nonlinear, wide-range IPG strips which give 2-D gel maps that are closely similar to the ones resulting with the CA-IEF technique used to establish the human keratinocyte database [15]. As an initial step towards interlaboratory comparisons of results obtained with the nonlinear gradient as a first dimension we report here on the focusing positions of 41 known proteins that are common to most human cell types. The pH range covered corresponds to the range in classical CA-IEF 2-D gel electrophoresis and in order to use these proteins as internal standards for comparing 2-D gel maps generated with other IPGs we determined their *pI* values with narrow-range IPGs in the first dimension. We have compared the calculated *versus* experimental *pI* values and show that it is necessary to have further information (absence or presence and nature of posttranslational modifications), in addition to amino acid composition to be able to calculate *pI* values that correspond to the actual experimental values. The *pK* values used for the calculations are provided and the usefulness of *pI* prediction in relation to database information is discussed. Furthermore, we comment on the possibility of using experimentally determined *pI* values to verify the available database information on polypeptide composition.

## 2 Materials and methods

### 2.1 Apparatus and chemicals

Equipment for isoelectric focusing and horizontal SDS electrophoresis (Multiphor<sup>®</sup> II electrophoresis chamber, Immobiline<sup>®</sup> strip tray, Multidrive XL programmable power supply, Macrodrive power supply and Multitemp<sup>®</sup> II) was from Pharmacia LKB Biotechnology AB (Uppsala, Sweden). Vertical second-dimensional gels were run in the home-made equipment described in [15]. The IPG strips with the wide-range nonlinear pH gradient were either Immobiline DryStrip<sup>®</sup> pH 3–10 NL, 180 mm or alternatively 160 mm long IPG strips with a corresponding pH gradient. In both cases the IPG strips were delivered by Pharmacia LKB Immobiline, Pharmalyte, Ampholine, GelBond as well as PAG film and the ready-made horizontal SDS gels (ExcelGel<sup>®</sup> XL SDS 12–14) were also from Pharmacia LKB. Purified proteins and peptides were from Sigma (St. Louis, MO).

### 2.2 Sample preparation

Preparation and labeling of unfractionated keratinocytes as well as fibroblasts have been described in [16]. Cells were lysed in a solution containing 9.8 M urea, 2% v/v NP-40, 100 mM DTT and 2% v/v Ampholine pH 7–9.

### 2.3 2-D gel electrophoresis

First-dimensional focusing was performed according to Gorg *et al.* [12] with some minor modifications, as described in [9]. Rehydration of the IPG strips was made in a solution containing 9.8 M urea, 2% v/v CHAPS, 10 mM DTT and 2% v/v carrier ampholyte mixture. The carrier ampholyte mixture consisted of 2 parts Pharmalyte

4-6.5. 1 part Ampholine pH 6-8 and 1 part Pharmalyte pH 8-10.5. Usually, cathodic sample application was used and the samples were diluted 2-20 times in a solution containing 9.8 M urea, 4% w/v CHAPS, 1% w/v DTT and 35 mM Tris base. For acidic application, the Tris-base was substituted with 100 mM acetic acid. The degree of dilution and sample volume (20-100  $\mu$ L) depended on the particular sample and the IPG, and whether visualization of the proteins was to be done by Coomassie Brilliant Blue or silver staining. With the wide-range non-linear IPG, 10-30  $\mu$ g of total protein was loaded for silver staining and 100-200  $\mu$ g for Coomassie staining. Focusing was done overnight with Vh products in the range of 45-60 kVh with 160 mm long strips and 50-70 kVh with 180 mm long strips. Solubilization of polypeptides and blocking of -SH groups prior to the second-dimensional run, as well as loading on the second-dimensional gel was done as described in [9]. The stacking gel was omitted and 5-10 mm were left at the top of the second-dimensional gel for applying the IPG strip. The space was filled with electrode buffer containing 0.5% w/v agarose. Casting, running, staining and autoradiography were carried out as described in [15].

## 2.4 Experimental determination of *pI* values

The determination of the *pI* differences between Immobilines *pI* 4.6, *pI* 6.2 and *pI* 7.0 necessary for the calibration of the *pH* scale at 25°C in 9.8 M urea was done as described in [9] with the same narrow-range IPGs. The *pH* scale was defined by setting the *pI* value of Immobililine *pI* 4.6 equal to 4.61 [9] and the determined *pI* differences gave the *pI* values of Immobilines *pI* 6.2 and *pI* 7.0, equal to 5.73 and 6.54, respectively. The *pI* differences found are in good agreement with values derived from [17] and [8] by extrapolation to 9.8 M urea concentration. As in [9], additional narrow-range recipes have been used for determining *pI* values. With narrow-range IPGs extending to *pH* values higher than the *pI* value of Immobililine *pI* 7.0, anodic sample application was used with acetic acid added to the sample solution. Otherwise, cathodic sample application was used with the same sample buffer as for wide-range IPGs.

## 2.5 Protein compositions used for *pI* calculations

With the exception of vimentin, protein compositions are from the Swiss-Prot database [18]. For vimentin, we used the data from [19], where the amino acid at position 41 is a D instead of a S. Information in the Swiss-Prot database on phosphorylation has been disregarded because it was known from earlier studies (J. E. Celis, unpublished results) that the spots in question corresponded to the unphosphorylated forms of the peptides.

## 2.6 Calculation of *pI* values

For the *pI* calculations it was assumed that the same *pI* value could be used for an amino acid residue in all polypeptides and in all positions in the peptide except for N- or C-terminally placed amino acids. For the *pI* values of the N-terminal amino groups the effect of the

different substituents on the  $\alpha$ -carbon were taken into account. The calculations of *pI* values were made with the aid of the IPG-maker program [20].

## 2.7 *pK* values used for *pI* calculations

For the carboxyl terminal group and internal glutamyl and aspartyl residues the same *pK* values were used as in [9]. For C-terminal glutamyl and aspartyl residues, separate *pK* values were derived with the aid of the Taft equations [9, 21]. The *pK* values of histidyl groups were calculated from the *pI* values of human carbonic anhydrase I as in [9]. For N-terminal glycine a *pK* value of 7.50 was used. The *pK* shift caused by a substituent on the  $\alpha$ -carbon was assumed to be identical with the *pK* shift the substituent caused for the amino group in the amino acid, i.e. 2.28 *pH* units were subtracted from the *pK* values for the amino groups in the amino acids given in [22, 23]. The approximate *pK* value of 9 for the cystenyl group was taken from [24]. For tyrosyl and arginyl groups we used the *pK* values for the amino acids [22, 23]. For lysyl groups the effect of high urea concentration on amino groups was taken into account and 0.5 *pH* units were subtracted from the amino acid *pK* value. These last three *pK* values are far from the *pH* range under study and the results found would have been the same if lysyl and arginyl groups were assumed to be fully ionized while the ionization of tyrosyl groups were neglected. A complete list of the *pK* values used is given in Table 1.

Table 1. *pK* values used for the ionizable groups in peptides 9.8 M urea, 25°C

Ionizable group	<i>pK</i>
C-terminal	
N-terminal	3.55
Ala	2.34
Met	2.00
Ser	6.03
Pro	8.36
Thr	6.82
Val	2.44
Glu	2.70
Internal	
Asp	4.05
Glu	4.45
His	5.98
Cys	9
Tyr	9
Lys	10
Arg	12
C-terminal side chain groups	
Asp	4.55
Glu	4.75

## 2.8 Statistical analysis

Statistical comparisons of the experimental and calculated *pI* values were done on an Apple Macintosh IIx using the statistical package Statistica/Mac, release 3.0b (from StatSoft Inc., Tulsa, Oklahoma). Calculated and experimental *pI* values were compared by the *t*-test for

correlated samples (paired *t*-test). The normality of *pI* differences was estimated graphically by probability plots. The variances of the data presented here and the similar data on plasma and liver proteins in [9] were compared by the *F*-test.

### 3 Results and discussion

#### 3.1 Identification of polypeptides and *pI* determinations

The 2-D gel maps of [<sup>35</sup>S]methionine-labeled proteins from noncultured, unfractionated normal human kerati-

nocytes, focused with the nonlinear, wide-range IPG and CA-IEF pH gradients in the first dimension, are shown in Figs. 1 and 2, respectively. The IPG extends to higher *pH* values but otherwise the two patterns are very similar and most of the spots in the IPG pattern can be directly related to the corresponding patterns in the CA-IEF gel. To obtain comparable patterns it was important to keep the focusing temperature as similar as possible. Compared to other studies [1-4, 9, 10, 12-14], we increased the urea concentration in the focusing gel to 9.8 M because keratins streaked badly in the focusing dimension when 8 M urea was used, presumably due to

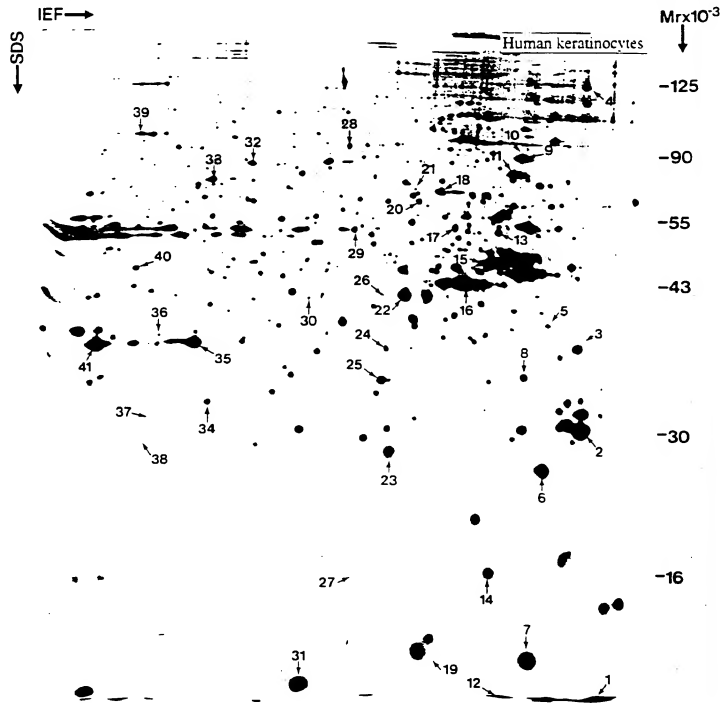


Figure 1. 2-D gel protein map of [<sup>35</sup>S]methionine-labeled proteins from noncultured, unfractionated normal human keratinocytes focused with the nonlinear, wide-range IPG in the first dimension. The position of the 41 proteins analyzed in this study is indicated.

aggregates of acidic and basic keratins. An increase in urea concentration to 9 M or more eliminated these streaks; apart from this effect, no other major changes in the focusing positions were observed. In Fig. 1 we have indicated the positions of 41 known proteins from the human keratinocyte 2-D gel database that are most likely common to most human cell types. The choice was made because these proteins are easy to identify with certainty. With the exception of stratifin (spot 2), involucrin (spot 4) and keratin 14 (spot 15), which are all

epithelial markers, these proteins are also present in human fibroblasts (Fig. 3) and lymphocytes (results not shown), and therefore can be used as landmarks for comparing 2-D gel maps derived from different cell types. In Table 2 the 41 proteins are listed together with their sample spot numbers (SSP) in the human keratinocyte protein database and *pI* values determined in 2-D gel maps generated with narrow-range IPGs in the first dimension.

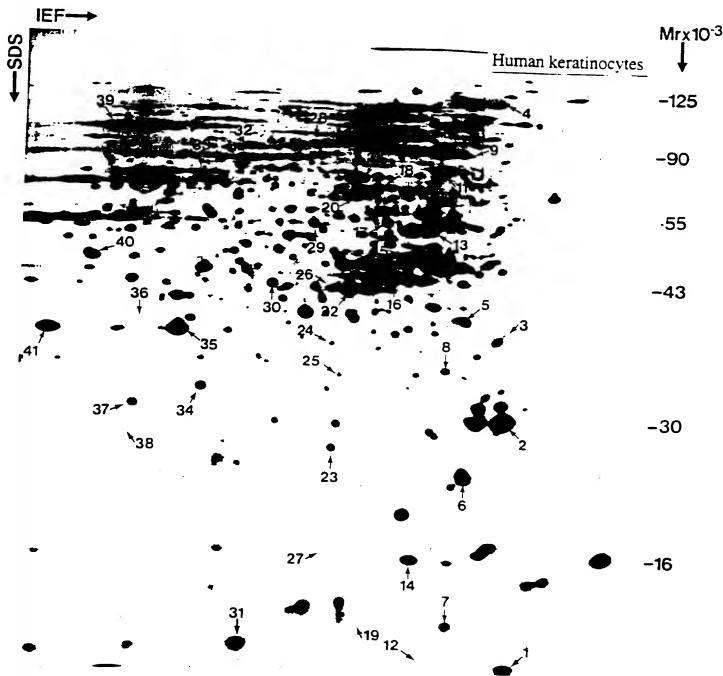


Figure 2. 2-D gel protein map of [ $^{35}$ S]methionine-labeled proteins from noncultured, unfractionated normal human keratinocytes focused with CA-IEF in the first dimension. The position of the 41 proteins analyzed in this study is indicated.

Table 2. Proteins from the human keratinocyte database localized in 2-D gels run with IPGs, as first dimension. Number in Protein name Figs. 1-3

Number	Protein name	H1 XSP number	Experimental pI value	Calculated pI value	Discrepancy (pI units)	Calculated net charge at experimental pI value	Inlet capacity charge units per pI unit	N-terminal	Recalculated for suspected blockage	N-terminal	Discrepancy per charge pI units	Net charge	Swiss Prot accession number
1	Cap 19	9027	4.46										
2	Caprin, type II (L33) related protein	9109	4.58										
3	Protein kinase alpha (PKA) (cyclic)	9226	4.58	4.57	-0.01	-0.1	20.8	M					P12004
4	Protein kinase alpha (PKA) (cyclic)	9226	4.58	4.63	0.05	0.3	20.1	M					P12004
5	Nucleolar protein B23	8211	4.75	4.64	-0.11	-1.2	30.4	M					P12004
6	Randomly coiled linear protein	8006	4.76	4.84	0.05	0.6	13.1	M <sup>2</sup>					P12004
7	Thiostatin	8213	4.85	4.84	-0.04	-0.1	20.3	M <sup>2</sup>					P12004
8	Annexin V	8611	4.95	4.88	-0.07	-0.3	20.3	M <sup>2</sup>					P12004
9	Heat shock protein 90-β	8213	4.95	4.91	-0.04	-0.1	20.3	M <sup>2</sup>					P12004
10	Heat shock protein 90-α	7629	4.97	4.97	0.00	0.0	36.2	M <sup>2</sup>					P12004
11	Glucose regulated protein 78 (Hsp70)	8315	4.99	4.98	-0.01	-0.1	31.6	M <sup>2</sup>					P12004
12	Calycin	8017	5.02	5.12	0.10	0.6	13.1	M <sup>2</sup>	5.09	0.07	0.3		P12004
13	Vimentin	8117	5.05	5.06	0.01	0.2	27.1	M <sup>2</sup>					P12004
14	Inhibition factor 1D	8006	5.08	5.08	0.00	0.0	27.1	M <sup>2</sup>					P12004
15	Keratin 14	7105	5.08	5.09	0.01	0.2	21.0	M <sup>2</sup>					P12004
16	β-Actin	7316	5.21	5.21	0.00	0.06	13.3	M <sup>2</sup>					P12004
17	Heat shock protein 60	6301	5.21	5.21	0.00	0.1	17.5	M <sup>2</sup>					P12004
18	Heat shock protein 70 (Hsp70)	6301	5.28	5.17	-0.09	-1.8	18.1	M <sup>2</sup>					P12004
19	Heat shock protein 70 (Hsp70)	6301	5.30	5.18	-0.08	-0.2	3.0	M <sup>2</sup>	5.32	0.04	0.8		P12004
20	Protein	5612	5.31	5.11	-0.07	-1.3	17.7	M <sup>2</sup>	5.36	0.02	0.3		P12004
21	Calreticulin	5612	5.31	5.11	-0.07	-1.3	17.7	M <sup>2</sup>	5.36	0.02	0.3		P12004
22	Plasminogen activator inhibitor 2	6111	5.38	5.36	-0.02	-0.3	23.3	M <sup>2</sup>	5.37	0.01	0.07		P12004
23	Glutathione S-transferase 3	5101	5.41	5.36	-0.05	-0.9	10.7	M <sup>2</sup>	5.37	0.01	0.07		P12004
24	Annexin VIII	5211	5.41	5.36	-0.05	-0.9	10.7	M <sup>2</sup>	5.37	0.01	0.07		P12004
25	Annexin III	5201	5.46	5.45	-0.01	-0.1	8.7	M <sup>2</sup>	5.46	0.01	0.05		P12004
26	Adenosine deaminase	5105	5.47	5.61	0.16	1.8	10.8	M <sup>2</sup>	5.52	0.06	0.5		P12004
27	Stathmin	5001	5.55	5.61	0.06	0.4	16.5	M <sup>2</sup>	5.51	0.07	0.8		P12004
28	Gelsolin, cytoplasmic	5608	5.55	5.58	0.03	-0.1	16.5	M <sup>2</sup>					P12004
29	Bar phosphatase-specific protein binding	5111	5.62										P12004
30	Protein	5111	5.74										P12004
31	Protein	1006	5.75										P12004
32	Cyclin, G1	5101	5.99	5.95	-0.04	-0.5	13.2	M <sup>2</sup>					P12004
33	Protein	5115	6.11	6.09	-0.02	-0.2	9.8	M <sup>2</sup>					P12004
34	Protein	5115	6.11	6.09	-0.02	-0.2	9.8	M <sup>2</sup>					P12004
35	Annexin I	2106	6.11	6.15	0.04	1.8	1.1	M <sup>2</sup>	6.28	0.17	0.9		P12004
36	Adenosine deaminase	5105	6.18	6.04	-0.16	-1.6	2.5	M <sup>2</sup>	6.31	0.15	0.6		P12004
37	Phosphoglycerate kinase (H form)	1107	6.35	6.35	0.00	0.7	4.2	M <sup>2</sup>	6.36	0.01	0.2		P12004
38	Triosephosphate isomerase	1111	6.31	6.35	0.04	0.4	2.6	M <sup>2</sup>	6.36	0.01	0.2		P12004
39	Phosphatase factor 2	1610	6.31	6.38	0.05	0.4	2.3	M <sup>2</sup>	6.36	0.01	0.2		P12004
40	α-Enolase	1125	6.62	6.62	0.00	1.5	9.8	M <sup>2</sup>					P12004
41	Annexin II	210	7.30	7.36	0.06	0.05	2.2	M <sup>2</sup>	6.75	0.11	0.1		P12004

at SSP number in the Keratinocyte database [15] at Peptides N-terminally sequenced as liver proteins [11]. Peptides given as N-terminally blocked in Swiss-Prot database.

### 3.2 Comparison between the determined and calculated $pI$ values for human keratinocyte proteins

Thirty six of the 41 proteins listed in Table 2 are found in the Swiss-Prot database. Contrary to the plasma and liver proteins used in [9], the  $pI$  calculations on the proteins used in this study posed some problems that reflected the way in which they were characterized. The

proteins used by Bjellqvist *et al.* [9] were either very abundant and well-characterized plasma proteins or they were identified by  $N$ -terminal sequencing and, therefore, the nature of the  $N$ -terminals (acetylated or non-acetylated) was in both cases known. The proteins used in this study have all been characterized by internal sequencing [7] and it is known that  $N$ -terminal acetylation occurs with high frequency in eukaryotes.

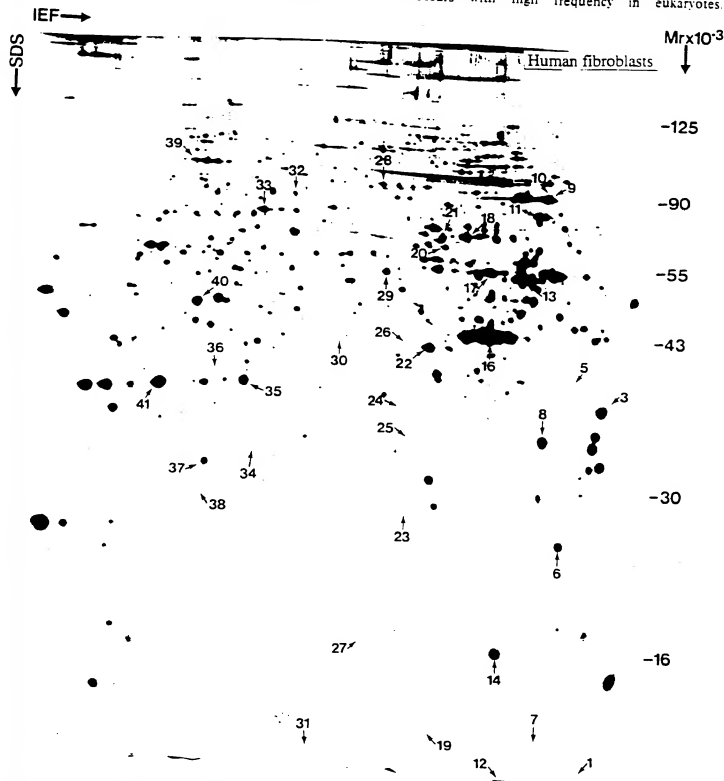


Figure 2 2-D protein map of [<sup>35</sup>S]methionine-labeled proteins from normal human fibroblasts focused with the nonlinear, wide-range IPG in the first dimension. The position of the 41 proteins analyzed in this study is indicated.



According to Brown and Robert [25], proteins with acetylated *N*-terminals correspond in weight to approximately 80% of the soluble protein in ascites cells. Based on results from *N*-terminal sequencing, at least 40% of the spots in the human liver protein 2-D gel map appear to be blocked [3]. The corresponding number, derived from 107 spots in the 2-D gel map of human T-lymphocyte proteins, falls between 60 and 65% (J. Strahler, personal communication). Information concerning *N*-terminal blockage is not normally available, and in the Swiss-Prot database only 6 of the 36 keratinocyte proteins are specified as *N*-terminally blocked. We have, within the present material, defined 18 proteins for which the *N*-terminals are very likely to be correctly described. Six of these proteins are listed in the Swiss-Prot database as *N*-terminally blocked, four represent proteins which appear in the human liver 2-D gel map and have been *N*-terminally sequenced as liver proteins [3] and the remaining eight have *N*-terminal groups other than M, S and A, *i.e.* *N*-terminals for which *N*-acetylation is uncommon [26]. In Figs. 4A, B, C and D *pI* values calculated from Swiss Prot database information are plotted against the experi-

mentally determined *pI* values for all the keratinocyte proteins listed in Table 2 and for the 18 selected proteins, as well as for the plasma and liver proteins (data from [9] valid for 10°C)\*.

The calculations show that without knowledge of the status of the *N*-terminal group, precise predictions of *pI* values for eukaryotic proteins cannot be achieved based on the information available in Swiss-Prot and similar databases. However, for proteins where the *N*-terminal status is known, we find good correlation between predicted and experimental *pI* values. When the variance of the *pI* discrepancies and the variance of calculated charges at the experimental *pI* values derived from the present data set are compared with the corresponding

\* There are four plots: (A) the 36 polypeptides from normal human keratinocytes (no corrections), (B) the 36 polypeptides from Fig. 4A where *pI* values have been recalculated for 12 polypeptides with M, S and A as *N*-terminally assumed blocked, based on calculated charge, (C) the 18 selected polypeptides with information on the *N*-terminal configuration, and (D) plasma and liver proteins.

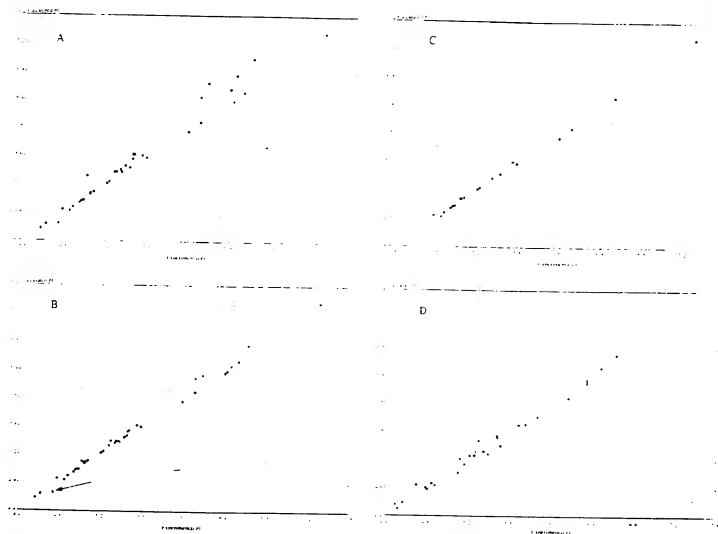


Figure 4. Calculated vs. experimental *pI* values. Lines are fitted using the least squares criterion: (A) 36 polypeptides from normal human keratinocytes (no corrections), (B) 36 polypeptides from Fig. 4A (including the 18 marker polypeptides where *pI* values have been recalculated assuming *N*-terminal blockage; x indicates recalculated *pI* values; nucleolar protein B23 is indicated with an arrow), (C) 18 polypeptides with information on *N*-terminal configuration, and (D) plasma and liver proteins.



composition used in the calculation is correct and complete. Exceptions to this are proteins such as involucrin and heat shock protein 90 that have very high buffer capacities. Introduction of an extra charge unit into these proteins will only result in *pI* shifts falling in the range of 0.01–0.02 pH units and the effect is that the quality of the *pH* definition – the precision by which *pK* values used in the calculations are given and the precision of experimental *pI* values in these cases – will limit the possibilities to verify polypeptide composition based on the experimental *pI* value.

Statistical comparison of experimental and calculated *pI* values was done using the *t*-test for dependent samples and normality of the discrepancies was estimated by probability plots. For the 36 proteins, the *p*-level is 0.0021, indicating that a result like this is unlikely to be a chance effect and must be assumed to represent a real difference. After correction for the most likely *N*-terminal configuration, the *p*-level is 0.043 and cannot be accepted as representing the same population since the *p*-level is less than 0.05 – the traditional *p*-limit of statistical significance. For the 18 proteins with a known or very likely *N*-terminal configuration the *t*-test gave a *p*-level of 0.49, which verifies that the experimental and calculated *pI* values are not significantly different.

Besides showing that *pI* values for denatured proteins with known compositions can be calculated with a high degree of precision from average *pK* values, the results also provide strong support for the notion that *N*-terminal blockage heavily depends on the nature of the *N*-terminal groups [26]. The results seem to indicate that with *N*-terminals other than M, S and A, only a few proteins have blocked *N*-terminals (1 out of 10 proteins in the present study), while it can be inferred from the data presented in Table 2 that a majority of the proteins with M, S and A as *N*-terminal are blocked. After correction for the effect of suspected *N*-terminal blockage there is only one protein (nucleolar protein B23) out of the 36 used in this study, which, in spite of a high buffer capacity, has a marked difference of 0.11 pH units between predicted and determined *pI* values (Fig. 4B), this corresponds to 3 charge units due to the high buffer capacity of this protein. This discrepancy in *pI* prediction and calculation of net charge at the *pI* is probably not due to deficiencies in the database information but instead reflects a shortcoming of the model used for *pI* calculations. Nucleolar protein B23 contains a domain extremely rich in aspartic and glutamic acid residues (Table 4), in which 26 out of 28 amino acid residues from position 161 to 188 are either a D or an E. A calculation based on the use of average *pK* values uninfluenced by the charged neighboring amino acid residues cannot be expected to correctly describe the *pI* value with almost half of the acidic groups packed

together into a highly negatively charged region. This limitation caused by calculations based on average *pK* values does not severely limit the usefulness of the approach since a search through Swiss-Prot shows that this type of D/E-rich motif is uncommon, and the existence of a highly charged region is immediately apparent upon inspection of the amino acid sequence.

The quality of the information available in databases, especially concerning posttranslational modifications, is a major problem when the data is to be used for *pI* predictions. The *p*-level of 0.043 found for all 36 proteins after correction for *N*-acetylation, shows that this problem is not only limited to *N*-terminal blockage and the very good agreement found for the eighteen polypeptides, with assumingly correctly described *N*-terminal (Fig. 4C), must be regarded as an exception from this point of view. *N*-terminal blockage is generally the main problem in relation to *pI* predictions for eukaryotic proteins. Of the 36 keratinocyte proteins analyzed, 18–20 are suspected to be *N*-terminally blocked (6 proteins blocked according to Swiss-Prot, 12 proteins with M, S or A as *N*-terminal and assumingly blocked based on the calculated charge, and two proteins, involucrin and nucleolar protein B23, with M as *N*-terminal for which the data does not allow any conclusion). This is in reasonable agreement with the conclusions based on the *N*-terminal sequencing data derived in connection with 2-D gel electrophoresis. *N*-terminal blockage can be suspected for 17–19 of the 26 proteins with M, S or A as *N*-terminal, while only 1 in 10 proteins with other *N*-terminal groups are blocked. The information that the frequency of *N*-terminal blockage is strongly related to the nature of the *N*-terminal group will be of some help in connection with *pI* predictions based on database information. However, without information from other sources, an uncertainty will always remain as to whether the *N*-terminal charge should be included in the *pI* calculation.

#### 4 Concluding remarks

The data presented here lays the foundation for comparing 2-D gel protein maps of different cell types generated with nonlinear, wide-range IPGs in the first dimension. The focusing positions of 41 polypeptides common to most human cell types have been described in a *pH* scale that allows focusing positions to be predicted with a high degree of accuracy, provided that the composition of the polypeptides are known and that information on posttranslational modifications are available. For polypeptides with a very high buffer capacity, the limiting factor is the precision with which experimental *pH* values can be determined rather than the precision of the calculations. Possible deficiencies in the *pH* scale description of the variation of the hydrogen ion activity has, at least at the present state, no consequences for its practical use. The major limitation in connection with predictions of focusing positions from polypeptide compositions is the quality of existing data on protein compositions, especially concerning posttranslational modifications. Amino acid sequences have been reasonably easy to obtain, while posttranslational modifications

Table 4. Amino acid sequence of nucleolar phosphoprotein B23

1	MDKDDG	APPNDP	DEADN	PDDEE	DEADN
51	ADGENT	EDAVES	PDGAG	NDPDE	DEDFG
101	PDGAGP	PDGAGP	PDGAGP	PDGAGP	PDGAGP
151	PDGAGP	PDGAGP	PDGAGP	PDGAGP	PDGAGP
201	PDGAGP	PDGAGP	PDGAGP	PDGAGP	PDGAGP
251	PDGAGP	PDGAGP	PDGAGP	PDGAGP	PDGAGP

have been difficult and work-intensive to determine. Recent developments in the field of mass spectrometry are fast changing this situation and within the next years we can expect a surge in reliable data in this area. While awaiting this development, verification of correctness and completeness of available information on polypeptide composition can be provided by experimental  $pI$  values in a pH scale based on the  $pI$  values determined in this study. So far, our data cover the pH range below  $pH = 7.5$ . The basic pH range covered by NEPHGE as first dimension will be covered in forthcoming work.

Received December 29, 1993

## 5 References

- [1] Gianazza, E., Astrua-Testori, S., Caccia, P., Giacomini, P., Quagliari, L., Righetti, P. G., *Electrophoresis* 1986, 7, 76-83.
- [2] Gorg, A., Postel, W., Gunther, S., *Electrophoresis* 1988, 9, 531-546.
- [3] Hochstrasser, D. F., Frutiger, S., Pauet, N., Barroch, A., Ravier, F., Pasquali, C., Sanchez, J.-C., Tissot, J.-D., Bjellqvist, B., Vargus, R., Appel, R. D., Hughes, G. J., *Electrophoresis* 1992, 13, 992-1001.
- [4] *Immobilizing Dr-Strip Kit for 2-D Electrophoresis*, Instructions, Pharmacia LKB Biotechnology AB, Lysala 1993.
- [5] Anderson, N. L., Hickman, B. J., *Anal. Biochem.* 1979, 91, 312-320.
- [6] Neidhardt, F. C., Appleby, D. A., Sunkar, P., Hutton, M. E., Phillips, T. A., *Electrophoresis* 1989, 10, 116-121.
- [7] Rasmussen, H. H., Damme, J. V., Puyse, M., Gesser, B., Celis, J. E., Vandekerckhove, J., *Electrophoresis* 1992, 13, 966-969.
- [8] Gianazza, E., Arioni, G., Righetti, P. G., *Electrophoresis* 1983, 4, 321-326.
- [9] Bjellqvist, B., Hughes, G. J., Pasquali, C., Pauet, N., Ravier, F., Sanchez, J.-C., Frutiger, S., Hochstrasser, D. F., *Electrophoresis* 1993, 14, 1023-1031.
- [10] Bjellqvist, B., Pasquali, C., Ravier, C., Sanchez, J.-C., Hochstrasser, D. F., *Electrophoresis* 1993, 14, 1357-1365.
- [11] O'Farrell, P. H., *J. Biol. Chem.* 1975, 250, 4007-4021.
- [12] Gorg, A., *Biochem. Soc. Transactions* 1993, 21, 135-132.
- [13] Hanash, S. M., Strahler, J. R., Neel, J. V., Hallat, N., Mathem, R., Keim, D., Zhu, X. X., Wagner, D., Gage, D. A., Watson, J. T., *Proc. Natl. Acad. Sci. USA* 1991, 88, 5709-5713.
- [14] Gorg, A., Postel, W., Friedrich, C., Knick, R., Strahler, J. R., Hanash, S. M., *Electrophoresis* 1991, 12, 653-655.
- [15] Celis, J. E., Rasmussen, H. H., Olsen, E., Maassen, P., Leffers, H., Honore, B., Dejgaard, K., Gromos, P., Hoffmann, H. J., Nielsen, M., Vassiles, A., Nintermyr, O., Hao, J., Celis, A., Basse, B., Lauridsen, J. B., Ratz, G. P., Andersen, A. H., Walbum, E., Klarsgaard, I., Puyse, M., Van Damme, J., Delay, B., Vandekerckhove, J., *Electrophoresis* 1993, 14, 1091-1105.
- [16] Celis, J. E., Maassen, P., Rasmussen, H. H., Leffers, H., Honore, B., Gesser, B., Dejgaard, K., Olsen, E., Magnusson, N., Knil, J., Celis, A., Lauridsen, J. B., Basse, B., Ratz, G. P., Andersen, A., Walbum, E., Brandstrup, B., Pedersen, P. S., Brandt, N. J., Puyse, M., Van Damme, J., Vandekerckhove, J., *Electrophoresis* 1991, 12, 802-872.
- [17] Bjellqvist, B., Ek, K., Righetti, P. G., Gianazza, E., Gorg, A., Postel, W., Westermeyer, R., *J. Biochem. Biophys. Methods* 1982, 6, 317-333.
- [18] Barroch, A., Boeckman, B., *Nucleic Acids Res.* 1991, 19, 2247-2249.
- [19] Honore, B., Maassen, P., Basse, B., Andersen, A., Walbum, E., Celis, J. E., Leffers, H., *Nucleic Acids Res.* 1990, 18, 6692.
- [20] Altland, K., *Electrophoresis* 1990, 11, 140-147.
- [21] Perrin, D. D., Dempsey, B., Serjant, E. P., *pH: Predictions for Organic Acids and Bases*, Chapman and Hall Ltd., London 1981.
- [22] Perrin, D. D., *Dissociation Constants of Organic Bases in Aqueous Solutions*, Butterworths, London 1965.
- [23] Perrin, D. D., *Dissociation Constants of Organic Bases in Aqueous Solutions*, Supplement 1972, Butterworths, London 1972.
- [24] Altland, K., Becker, P., Rossman, U., Bjellqvist, B., *Electrophoresis* 1988, 9, 474-485.
- [25] Brown, J. L., Robert, W. K., *J. Biol. Chem.* 1976, 251, 1009-1014.
- [26] Persson, B., Flinta, C., Heine, G., Jorvall, H., *Eur. J. Biochem.* 1985, 152, 527-537.

Bo Franzén<sup>1</sup>  
 Stig Linder<sup>2</sup>  
 Ken Okuzawa<sup>2</sup>  
 Harabumi Kato<sup>2</sup>  
 Gert Auer<sup>1</sup>

<sup>1</sup>Division of Tumor Pathology,  
 Department of Pathology, Division  
 of Experimental Oncology,  
 Karolinska Hospital and Institute,  
 Stockholm Sweden  
<sup>2</sup>Tokyo Medical College, Department  
 of Surgery, Tokyo  
<sup>3</sup>Division of Experimental Oncology,  
 Karolinska Hospital and Institute,  
 Stockholm

## Nonenzymatic extraction of cells from clinical tumor material for analysis of gene expression by two-dimensional polyacrylamide gel electrophoresis

We have compared different methods of preparation of malignant cells for two-dimensional electrophoresis (2-DE). We found all methods using fresh tissue to be superior compared to methods using frozen tissue. Our results indicate that nonenzymatic methods of preparation of tumor cells, including fine needle aspiration, scraping and squeezing, have advantages over methods using enzymatic extraction of cells. Nonenzymatic methods are rapid, appear to reduce loss of high molecular protein species, and alleviate the necessity of separating viable and nonviable cells by Percoll gradient centrifugation. Using these techniques, high-quality 2-DE maps were derived from tumors of the lung and breast. In the resulting polypeptide patterns, heat shock proteins, non-muscle tropomyosins and intermediate filament were identified. We conclude that nonenzymatic extraction of malignant cells from fresh tumor tissue improves the possibilities that these techniques may be useful in clinical diagnosis.

### 1 Introduction

Tumors may develop by a number of different mechanisms in any given cell type. At the time of diagnosis, tumors will have progressed along different pathways to various stages of malignancy. To provide a basis for individual therapy it is of importance to examine specific properties of the tumor cell population in each patient. A large number of different markers have been described in order to increase the diagnostic accuracy. It is likely that a combination of several markers is needed in the future in order to reflect different properties of the tumor. One important method for the resolution of a large number of potential markers is two-dimensional electrophoresis (2-DE). Extensive efforts are being made in identifying various polypeptides separated by 2-DE and to characterize how the expression of these polypeptides is affected by the response to cellular transformation and various culture conditions [1,2]. It would be of value to transfer this information to 2-DE separations of polypeptides from tumor tissue samples. However, one prerequisite is that the quality of the 2-DE gels from tumor samples is comparable in quality with 2-DE gels from samples of cultured cells.

Frozen tumor tissues are commonly used for various biochemical assessments. However, if such samples are analyzed by 2-D polyacrylamide gel electrophoresis (PAGE), the polypeptide patterns are obscured by contamination of serum- and connective tissue proteins. Such nontumor-cell-related variations represent serious problems in the interpretation and inter-patient comparison of 2-DE

patterns [3]. 2-DE patterns of cells prepared from fresh tumor material were analyzed after enzymatic extraction of tumor cells [4, 5] or after culturing tumor fragments in medium containing radioactive amino acids [6]. These procedures may, however, lead to alterations in the gene expression/polypeptide patterns. We are only aware of one study where nonenzymatic extraction of cells from fresh tumor tissue (prostate cancer) was used to prepare samples for 2-D PAGE [4]. We have examined enzymatic extraction and various nonenzymatic preparation techniques, including fine needle aspiration, for the preparation of cells from fresh tumor tissues. We describe nonenzymatic extraction procedures that are rapid, lead to high-quality 2-DE patterns, and that alleviate the necessity to purify tumor cell populations from dead cells.

### 2 Materials and methods

#### 2.1 Cell cultures and samples used for spot identification

A rat embryonal fibroblast cell line, WT2 (a kind gift from Dr. J. I. Garrels and Dr. S. Patterson) was used for the identification of a number of heat shock and structural proteins. Human normal diploid lung fibroblasts, WI38, human epithelial breast carcinoma cells, MDA-231 and MCF-7 were purchased from ATCC and grown as recommended. Polypeptides prepared from a leukemia type pre-B-ALL were separated by 2-DE. The 2-DE map was then analyzed by Dr. S. M. Hanash (University of Michigan, Ann Arbor, USA).

#### 2.2 Tumor tissues samples

In this study, 2-DE maps from seven tumors were used as representative illustrations: two adenocarcinoma of the lung (LA and LB, mucinous, both cases intermediate grade of differentiation), one squamous carcinoma of the lung (LS), one carcinoma-like breast cancer (BC), one microfollicular adenoma (highly differentiated) of the thyroid (TA), one highly differentiated hyperneph-

Correspondence: Dr. Bo Franzén, Division of Tumor Pathology, Department of Pathology, L1:01, Karolinska Hospital and Institute, 10401 Stockholm 60, Sweden

Abbreviations: 2-DE, Two-dimensional polyacrylamide gel electrophoresis; IEF, isoelectric focusing; LDH, lactate dehydrogenase; NP-40, Nonidet P-40; PBS, phosphate buffered saline; PCNA, proliferating cell nuclear antigen; PLH, protease inhibitors; PMSF, phenylmethyl sulfonyl fluoride; SDS, sodium dodecyl sulfate; WW, wet weight

© VCH Verlagsgesellschaft mbH, 0945 Weinheim, 1993

0173-0825/93/1010-1045 \$5.00+35.00

roma, a tumor of the kidney (KH), and finally one case of poorly differentiated corpus carcinoma (CP).

### 2.3 Preparation of cultured cells

The cell monolayers were washed twice in phosphate buffered saline (PBS) and then scraped off in ice-cold PBS including protease inhibitors (PIH), phenylmethylsulfonyl fluoride (PMSF) 0.2 mM and 0.83 mM benzamide pelleted at  $660 \times g$  for 3 min ( $+4^\circ\text{C}$ ) and washed one time before final centrifugation at  $2700 \times g$  for 5 min. The wet weight of the cell pellet was recorded and the cells were stored at  $-80^\circ\text{C}$  until further processing.

### 2.4 Preparation of tumor tissue samples

#### 2.4.1 General remarks

Macroscopically representative and non-necrotic tumor tissues were selected within 20 min after resection. Parallel samples were routinely prepared for cytology. The samples were processed as rapidly as possible on ice or at  $+4^\circ\text{C}$  and in the presence of PIH. Cells were stained with DiffQuick (Baxter) and usually examined at three different occasions during the preparation procedure: (i) cytology sample, (ii) extracted cells and (iii) cells after percoll gradient centrifugation.

#### 2.4.2 Specimen acquisition

The strategy of sample preparation is shown in Fig. 1. Tumor tissue cell samples were usually obtained by fine needle aspiration (NA) using a 0.7 mm needle. The syringe was filled with 1–2 mL of ice-cold culture medium/PIH. We found that if a tumor appeared to be very fibrous it is difficult to extract enough cells for 2-DE analysis. In these cases, two alternative techniques were examined. (i) The tumor was cut in the middle and the fresh surface scraped (SC) by a scalpel. The cell-rich material was then transferred to ice-cold culture medium (L15 with 5% fetal calf serum)/PIH. (ii) A part of the tumor sample was placed in culture medium on ice for further processing at the laboratory in the following way: the material was cut into very small fragments on a pre-cooled dissection plate and transferred to a small glass chamber with a 0.7 mm metal net 5 mm above the bottom of the chamber. Medium/PIH was added to cover the sample (8 mL) which was gently squeezed (SQ) towards the net in order to release and wash out cells. NA and SC were also compared with an enzymatic extraction (EE) procedure described previously [5]. Briefly, thin slices of tissue were incubated with collagenase (1 mg/mL) and elastase (2 mg/mL) in medium for 1 h at  $37^\circ\text{C}$ . Extracted cells from every sample were then subjected to percoll gradient centrifugation (Section 3.2.3).

#### 2.4.3 Separation of cells by Percoll gradient centrifugation

The cell suspension was filtered through two nylon mesh filters, (i) 250  $\mu\text{m}$  and (ii) 100  $\mu\text{m}$  and then centrifuged

at  $660 \times g$  for 3 min. The cell pellet was resuspended carefully in medium, using a syringe and loaded onto a two-step discontinuous Percoll/PBS gradient: 20.4 (density = 1.03 g/mL) and 54.7% (density = 1.07 g/mL) and centrifuged at  $1000 \times g$  for 15 min. In this system, dead cells stay on the top, viable cells sediment to the interphase and erythrocytes sediment to the bottom. The viability of cells in the top fraction and interphase was checked by the trypan blue exclusion test. The interphase cell layer ( $> 90\%$  viability) was collected and washed one time in a large volume PBS/PIH (centrifuged at  $800 \times g$  for 5 min). Finally, the cells were resuspended in 1.4 mL PBS and pelleted at  $2700 \times g$  for 5 min. The wet weight (WW) was recorded and the pellet was then stored at  $-80^\circ\text{C}$ .

#### 2.4.4 Final preparation of cells for 2-D PAGE analysis

From this point, cultured cell samples were treated in the same way as tumor cell samples. Each cell pellet was thawed on ice and resuspended in 1.89  $\mu\text{L}$  water per mg WW ( $= 1.89 \times \text{WW}$ )  $\mu\text{L}$ . The suspension was frozen and thawed 4–5 times to break the cells [7]. A volume of  $(0.089 \times \text{WW}) \mu\text{L}$  10% sodium dodecyl sulfate (SDS), including 33.3% mercaptoethanol, was mixed with the sample and incubated 5 min on ice with  $(0.329 \times \text{WW}) \mu\text{L}$  of a solution of DNase I (0.144 mg/mL 20 mM Tris-HCl with 2 mM  $\text{CaCl}_2 \times 2\text{H}_2\text{O}$ , pH 8.8) and RNase A (0.0718 mg/mL Tris) [8,9]. The sample was frozen and lyophilized. Sample buffer [10] including

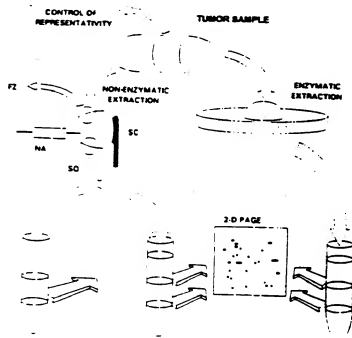


Figure 1. Experimental flow chart showing main steps of the preparation procedures. The abbreviations used for nonenzymatic extraction procedures are: FZ, frozen sample preparation; NA, needle aspiration; SC, scraped; and SQ, squeezed sample. Extracted cells are then loaded as a suspension (top volume of each tube) onto either 1.07 g/mL Percoll (left), or a discontinuous Percoll gradient from the nonenzymatic extraction (middle), or from enzymatic extraction (right). Cellular top- and interphase fractions are then used for 2-DE. For details see Section 2.

PMSF (0.2 mM, EDTA (1.0 mM), 0.5% Nonidet P-40 (NP-40), and 3-[3-cholamido propyl)-dimethylammonio]-1-propane sulfonate (CHAPS; 25 mM) was added carefully, mixed for 2.5 h and centrifuged for 15 min at

10000 rpm to remove any insoluble material. Duplicate or triplicate samples were taken for protein determination [11]. Samples were stored at  $-80^{\circ}\text{C}$  prior to isoelectric focusing (IEF).



Figure 2. 2-DE analysis of samples from three cell lines and one leukemia used for the identification of polypeptides: (A) WT2; (B) NDA-231. arrowheads mark some low molecular weight cytosolic polypeptides. (C) W38 and (D) pre B-ALL. The abbreviations for identified spots are explained in Table 1.

### 2.4.5 Preparation of frozen tumor tissue

The technique has been described previously [3,12]. Briefly, the sample is moarsted frozen to a fine powder, homogenized, lyophilized and solubilized in sample buffer.

### 2.4.6 Control of representativity

The tumors were examined routinely by experienced pathologists and smears or imprints from the samples were also assessed for cytometric DNA content by microspectrophotometry.

### 2.5 2-D PAGE

2-D PAGE was performed as described [8,10] except for the following details. The glass tubes for IEF, 1.2 × 200 mm, contained 2.0% Resolyte, pH 4–8 (BDH) and were cast to a height of 180 mm. A stock solution of acrylamide (Serva) and  $N,N'$ -methylenebisacrylamide (16.7:1 for IEF and 37.5:1 for the second dimension) was deionized by mixing with 5% w/v Duolite MB 5313 mixed-resin ion exchanger (BDH) for 30 min, filtered (with a 0.22  $\mu$ m nitrocellulose filter) and stored at  $-70^{\circ}\text{C}$ .  $N,N'$ -Methylenebisacrylamide,  $N,N,N',N'$ -tetramethylethylenediamine (TEMED) and ammonium persulfate were purchased from Bio-Rad. IEF tubes were prefocused at 200 V in 60 min. To each tube a sample corresponding to 20–40  $\mu$ g protein was applied and focused for 14.5 h at 800 V and finally 1.0 h at 1000 V using a Protean II cell (Bio-Rad) and Model 1000/500 Power Supply (Bio-Rad). The tube gels were finally extruded into 1.25 mL equilibration buffer, containing 60 mM Tris, pH 6.8 (2% SDS, 100 mM dithiothreitol and 10% glycerol), frozen on dry ice and stored at  $-70^{\circ}\text{C}$ . The second dimension (1.0 × 180 × 90 mm) of the acrylamide concentration was 10%

T, and the gel contained 376 mM Tris, pH 8.8, and 0.1% SDS. IEF gels were applied on top of the slab gel, sealed with 0.5% agarose containing electrophoresis running buffer (60 mM Tris-base, 0.2 M glycine and 0.1% SDS), and electrophoresed with 10–11 mA per gel (constant current) at  $-10^{\circ}\text{C}$ . Six gels were run together in a Protean II xi 2-D Multi-Cell (Bio-Rad). Proteins were visualized by silver staining and photographed with the acidic side to the left [13,14].

### 2.6 Identification of polypeptides

Vimentin and vimentin-derived polypeptides were identified by extraction of an MDA-231 cell lysate with 0.6 M KCl/0.5% NP-40 [15]. Tropomyosins were extracted from MDA-231 and W138 cell lysates [16], and cytokeratin were extracted from MDA-231 and MCF-7 cell lysates [17]. The patterns were compared with published maps [19–21]. Proliferating cell nuclear antigen (PCNA) was identified by immunoblotting (PC10 mAb, Dako-patt) using a semidry system (Multiphor II Nova Blot, Pharmacia-LKB Biotechnology AB) and enhanced chemoluminescence (ECL) detection (Amersham).

### 3 Results

#### 3.1 2-DE of samples prepared from normal and tumorigenic cultured cells

The object of this study was to develop methods for preparation of 2-DE maps from human tumor tissue which have the same high resolution as those obtained from cultured cells. Shown in Fig. 2 are high resolution 2-DE gels prepared from cultured cells and one leukemia: SV40 transformed embryonal rat fibroblast WT2 (Fig. 2a); human MDA-231 breast carcinoma cells (Fig. 2b); human W138 fibroblasts (Fig. 2c) and human pre B-ALL

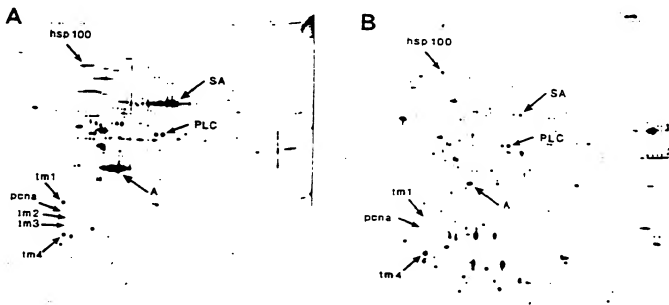


Figure 2. 2-DE analysis of a case of lung adenocarcinoma (LA). Comparison of 2-DE gel quality between (A) frozen and (B) fresh (needle aspiration) tissue preparation.



cells (Fig. 2d). Polypeptides were identified through a laboratory exchange of cell samples/2-DE maps and through 2-DE analysis of purified proteins (Table 1).

### 3.2 Preparation of samples from solid tumors

#### 3.2.1 Fresh versus frozen tissue

An adenocarcinoma of the lung (LA) was prepared for 2-DE by conventional methods using frozen material (Fig. 3a). There are several possibilities for the poor resolution using frozen tissue, including the presence of high molecular weight protein aggregates. Filtering extracts through 0.1  $\mu$ m filters (Durapore, Millipore) resulted in a slightly improved resolution (not shown). When fresh tumor tissue from tumor LA was used for sample preparation, using fine needle aspiration to collect the cells, the resolution was considerably improved (Fig. 3b). The use of fresh tissue resulted in a general increase in resolution, which was most pronounced in the 50–100 kDa molecular mass range. A number of differences in the protein profiles of the gels in Figs. 3a and 3b can be observed, some of which are indicated in the figures. The decrease in serum albumin in Fig. 3b is likely to result from loss of serum proteins occurring when cells were pelleted after aspiration. Other differences, such as the decreased level of transformation-sensitive tropomyosins (TM1-TM3), may result from enrichment of tumor cells in the sample of Fig. 3b. Fine needle aspiration, a well-established technique in cytology, extracts mainly tumor cells because of decreased intercellular adhesiveness of neoplastic cells as compared to normal tissue. Microscopic examination of Diff-Quick-stained extracted cells from case LA revealed almost 100% tumor cells, whereas the whole tissue extract contained approximately 60% tumor cells.

Table 1. Names and abbreviations for identified spots

Spot	Name	Basis for identification
A	Actins	a
aA	$\alpha$ -Actinin	a
B23	Protein B23/Nuclarn	a
EF2	Elongation factor 2	a
EF1	Elongation factor 1 $\beta$	a
GT	Glutathione-S-transferase ( $\mu$ )	a
hsp60	Heat shock protein 60	a
hsp73	Heat shock protein 73	a
hsp80	Heat shock protein 80, GRP78, BiP	a
hsp90	Heat shock protein 90	a
hsp100	Heat shock protein 100, Endoplasmic	a
IFa	Intermediate filament associated	a
k8	Cytokeratin 8	b and a
Lamb	Lamin B	a
Lip1	Lipocorin I	a
Lip2	Lipocorin II	a
Lip3	Lipocorin V	a
Mit1	Mitcon 1/8 – F1 ATPase	a
Mit2	Mitcon 2	a
Mit3	Mitcon 3	a
MRP	Mucine Related Polypeptides	a
pca	Proliferating cell nuclear antigen	c and a
PLC	Phospholipase C (1)	a
RO	RO/SS-A antigen	a
SA	Serum Albumin	b and a
aT	$\alpha$ -Tubulin	a
bT	$\beta$ -Tubulin	a
tm1	Non-muscle tropomyosin isoform 1	b and a
tm2	Non-muscle tropomyosin isoform 2	b and a
tm3	Non-muscle tropomyosin isoform 3	b and a
tm4	Non-muscle tropomyosin isoform 4	b and a
tm5	Non-muscle tropomyosin isoform 5	b and a
TPI	Triose phosphate isomerase	a
V	Vimentin	a
Vid1	Vimentin derived protein	b and a
Vid2	Vimentin derived protein	b and a
Vid3	Vimentin derived protein	b and a
Vid4	Vimentin derived protein	b and a
Vin	Vinculin	a

a, homologous position with respect to other mammalian systems

b, purified proteins

c, immunoblotting

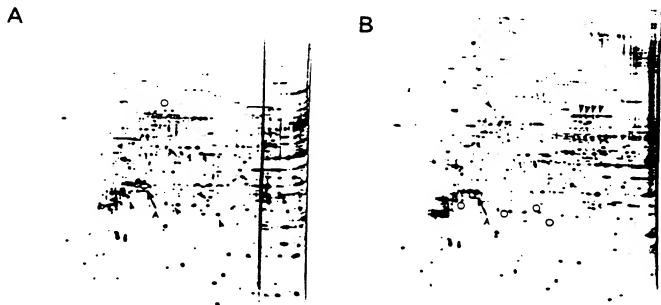


Figure 4. 2-DE analysis of a case of breast carcinoma (BC). Comparison of 2-DE quality and some differences in detected spots (arrow heads indicate increased intensity and circles or bracket indicate decreased intensity of the same spots) between (A) enzymatically and (B) nonenzymatically extracted tissue preparation.

### 3.2.2 Comparison of different methods for preparing cells from fresh tumor tissue

Samples were prepared from breast and lung carcinomas using either an enzymatic treatment with collagenase/elastase or using nonenzymatic preparations (Fig. 4). A number of differences in the protein profiles were observed in the resulting 2-DE gels, some of which are indicated in Figs. 4a and b. These differences include both increases and decreases in spot intensity. These differences may result from degradation of high molecular weight polypeptides during enzymatic treatment, increased solubilization of polypeptides, or may have other causes. For many tumors, it was only possible to obtain

small amounts of material since they were reserved for other examinations. In these cases, samples could be prepared for 2-DE using either needle aspiration or scraping. Figure 5a shows a 2-DE gel prepared from squamous lung carcinoma (LS) cells collected by needle aspiration and Fig. 5b shows a gel prepared from the same tumor by scraping. In this case, a number of differences were recorded between the two procedures, some of which are arrowed in Fig. 5. Samples obtained from other tumors (breast and lung) generally showed fewer differences between these two methods of cell sampling (not shown). These data show that different nonenzymatic extraction procedures may yield different polypeptide patterns. However, the number of spots with a large



Figure 5. 2-DE analysis of a case of lung cancer (LS). Comparison of 2-DE gel quality and detected spots (arrow heads and circles) between (A) aspirated (needle aspiration) and (B) scraped preparations from fresh tissue.

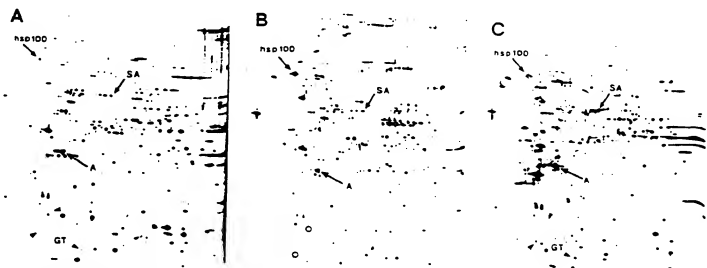


Figure 6. 2-DE analysis of three other types of tumors, (A) hypernephroma, (B) an adenoma of the thyroid and (C) corpus cancer, using the nonenzymatic preparation technique. Arrowheads and circles indicate some cytosolic polypeptides.

difference in intensity were lower than when a nonenzymatic preparation was compared with an enzymatic preparation.

2-DE maps of satisfactory quality were prepared by a third procedure. Cells were released from small pieces of tumor by squeezing (see Section 2). Some examples of this are shown in Fig. 6 where 2-DE maps derived from a case of hypernephroma, KH (Fig. 6a), a case of thyroid tumor, TA (Fig. 6b) and a case of corpus cancer, CP (Fig. 6c) can be seen. We conclude that nonenzymatic techniques are useful for 2-DE analysis of a number of different tumors. The quality of the resulting gels is com-

parable to that obtained using cultured cells (compare the gels in Fig. 2 with those in Fig. 4, 6 and 7). Which of these methods will be optimal will, in our experience, depend on the tumor material. For example, very small tumors are preferably extracted by squeezing; on the other hand, breast cancers (which are often fibrous) yield satisfactory samples using scraping.

### 3.2.3 Purification of cells on percoll gradients

We considered the possible advantage of separating viable cells from dead cells, erythrocytes, and debris using discontinuous Percoll gradients. Cells collected

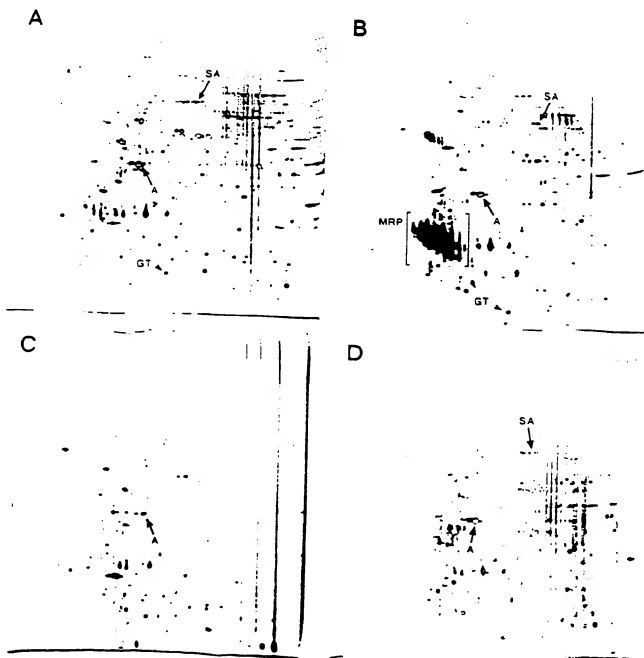


Figure 7 2-DE analysis of polypeptides from viable (b and d) and nonviable (a and c) cells of an adenocarcinoma of the lung (LB), separated using discontinuous Percoll density gradient. Nonenzymatic preparation technique (a and b) and enzymatic preparation technique (c and d) are compared.

from the interphase showed a viability of more than 90% as judged by trypan blue exclusion test. However, it was found that the yield of viable cells decreased dramatically if the tissue resection was not immediately processed. To study the effect of lysis of cells during the preparation procedure, 2-DE maps were prepared from nonenzymatically extracted cells of case LB collected from the top fraction (nonviable, Fig. 7a) and interphase fraction (viable, Fig. 7b). These 2-DE maps were compared with corresponding fractions (nonviable, Fig. 7c, and viable, Fig. 7d) of enzymatically extracted cells. One clear disadvantage of the enzymatic technique was that when loss of cell viability occurred during preparation, a dramatic loss of high molecular weight polypeptides was observed (Fig. 7c). This was probably due to degradation of intracellular proteins. However, nonenzymatic preparations showed fewer differences between viable and nonviable cells. The most pronounced alteration was a decrease of a group of mucine related proteins (Fig. 7b). We conclude, therefore, that discontinuous Percoll gradient is necessary after enzymatic extraction of cells, but can be omitted from the nonenzymatic tumor sample preparation procedure.

We used the MDA-231 cell line to study the effects of cell lysis and leakage of cytosolic polypeptides during sample preparation. Remarkably, after 30, 50, 80 and 140 min of incubation in PBS/PIH at 0°C, no significant changes were observed in the 2-DE pattern (not shown). Although loss of cell viability may not result in protein degradation when cells are incubated in the presence of protease inhibitors, loss of cytosolic proteins would be expected during pelleting of cells. We monitored the loss of lactate dehydrogenase (LDH) activity in the supernatant during incubation in PBS of MDA-231 and MCF-7 breast cancer cells at 20°C. In both cases, loss of viability was paralleled by release of LDH from the cells (Fig. 8). After 5 h, 70% of the MCF-7 cells, but only 30% of the MDA-231 cells were dead (not shown).

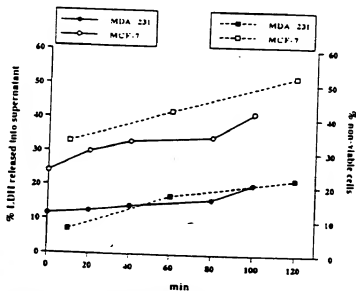


Figure 8. The relative release (fraction in supernatant) of total of lactate dehydrogenase activity (LDH) and cell viability versus incubation time of the mammary carcinoma cell lines MDA-231 and MCF-7 during incubation in PBS at 20°C.

These data indicate the impact of a rapid preparation procedure, at low temperature, of fresh tumor samples. Experiments have also been performed using only 1.07 g/mL Percoll (Fig. 6c and Fig. 1, left test tube) in order to remove erythrocytes. One clear advantage with this procedure, which today is routinely utilized, is a higher yield of viable cells, probably due to decreased sample preparation time.

#### 4 Discussion

We describe procedures for sample preparation from solid tumors for 2-DE. 2-DE maps could be derived from solid tumors which were similar in quality to those obtained from cultured cells. Compared to methods using frozen material, the resolving power of the 2-DE technique is increased, allowing examination of a large number of polypeptides from tumors of different malignancies. Other investigators [12,22] have used samples from frozen tumors to derive 2-DE maps. We have previously described disadvantages encountered using frozen tumor samples including variations in contaminating proteins between different samples [3]. The methods described here are based on the preparation of cells from tumors without enzymatic digestion. The enzymatic step could be avoided since malignant cells usually grow as solid masses which are not strongly attached to the matrix. Furthermore, we found that omitting the enzymatic digestion alleviated the necessity of purifying viable tumor cells on Percoll gradients. This was in sharp contrast to enzymatically treated samples, where loss of viability leads to loss of high molecular weight proteins (Fig. 7c).

At least in the case of lung cancer, viable and nonviable cells showed small differences in respect to 2-DE maps. Presumably, protease inhibitors penetrate cells and inhibit proteolysis. In model experiments, we observed leakage of cytosolic protein (LDH) from the cells in parallel to loss of viability. Apparently, however, only a limited decrease of the level of low molecular weight cytosolic polypeptides was detected using silver staining combined with visual inspection. We have found that although some tumors are well suited for the preparation procedure described, others are not. In general, good results were obtained using tumors of the lung, breast, corpus and lymphomas. In contrast, cells from thyroid adenomas and hypernephroma showed poor viability. We were in these cases unable to separate nonviable cells from viable cells, and we can therefore not evaluate the consequence of the loss of viability on 2-DE patterns, apart from a loss of some low molecular weight cytosolic polypeptides.

Highly differentiated tumors may show lower viability as compared with poorly differentiated tumors (Dr. Farkas Vanky, personal communication). A number of samples from thyroid tumors were prepared for 2-DE but most cases showed poor viability. We believe that special care is needed during preparation of generally highly differentiated tumor groups. The difference between loss of viability/leakage of LDH of the more differentiated MCF-7 cells and the less differentiated MDA-231 cells is in line

with these observations (Fig. 8). A number of potential and interesting markers, like tropomyosin isoforms, cytokeratins and heat shock proteins, appear to be insensitive to loss of viability during the preparation procedure. We have to date made numerous observations of alterations in the expression of these polypeptides in breast cancers and lung cancers.

Another problem that may occur, irrespective of sample preparation techniques used, is admixture of lymphocytes. These cases are easily detectable in smears and it may therefore be possible to select lymphocyte specific spots as "internal markers" for the 2-D PAGE analysis. Studies using this approach are in progress. Many of the polypeptides identified are structural (Table 1). Since the expression of many of these polypeptides are known to vary between normal and malignant cells, the possibility to determine their expression simultaneously is appealing. In the specific case of breast cancer, alterations in the expression of intermediate filament proteins (cytokeratins) are known to occur during tumor progression [23]. Other proteins known to be differentially expressed between normal cells and transformed cells are tropomyosins, numatrin/B23, heat shock proteins and PCNA. To this end, we have observed alterations in the expression of cytokeratin 8, hsp 90, and non-muscle tropomyosin isoform 2 during malignant progression. (Okuzawa *et al.*, in preparation and Franzen *et al.*, in preparation).

The method of choice for sample preparation from tumor tissues will depend on the properties of the tumor material studied. It may be important to use only one method when comparing cases within one group, as differences were observed between methods. The advantages of the nonenzymatic techniques are (i) that it minimizes contamination with connective tissue, (ii) that problems with contamination of serum proteins are avoided, and (iii) that separation of viable and dead cells is not necessary. Hereby the resolving power of 2-D PAGE is maximized for the analysis of human tumors and studies on inter-tumor variations in gene expression are facilitated. In addition, the polypeptide patterns obtained may be more representative for the *in vivo* tumor cell since the use of enzymes and incubations have been minimized.

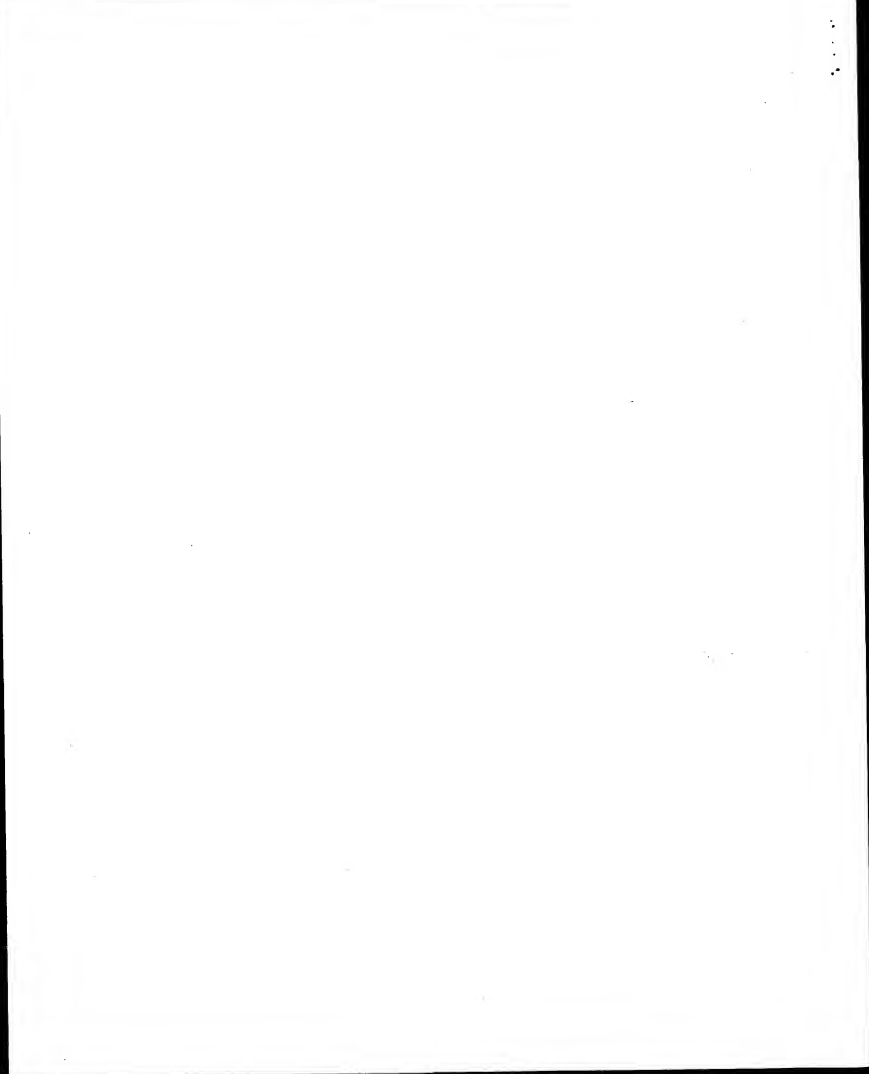
We would like to thank Dr. J. I. Garrels, Dr. S. Patterson, Dr. S. M. Hanash and Dr. J. E. Celis for making sample and 2-DE map exchanges possible. This study was sup-

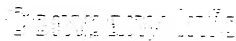
ported by grants from the Swedish Cancer Society and the Cancer Society in Stockholm.

Received March 5, 1993

## 5 References

- [1] Celis, J. E., Deigaard, K., Madsen, P., Leffers, H., Gesser, B., Honore, B., Rasmussen, H., Olsen, E., Lauridsen, J. B. and Ratz, G., *Electrophoresis* 1990, 11, 1072-1115.
- [2] Garrels, J. I., Franza, B. R., Chang, C., Litter, G., *Electrophoresis* 1990, 11, 1114-1130.
- [3] Franzen, B., Iwabuchi, H., Kato, H., Lindholm, J. and Auer G., *Electrophoresis* 1991, 12, 509-515.
- [4] Sherwood, E. R., Berg, L. A., Mitchell, N. J., McNeal, J. E., Kozlowski, J. M. and Lee, C. J., *Urology* 1990, 143, 167-171.
- [5] Endler, A. T., Young, D. S., Wold, L. E., Lieber, M. M. and Currie, R. M., *J. Clin. Chem. Clin. Biochem.* 1986, 24, 981-992.
- [6] Forchhammer, J. and Macdonald-Bravo, H., in: *Celis, J. E. and Bravo, R. (eds.), Gene Expression in Normal and Transformed Cells*, Plenum, New York 1983, pp. 291-314.
- [7] Linder, S., Brzeski, H. and Ringertz, N. R., *Exp. Cell Res.* 1979, 120, 1-14.
- [8] Celis, J. E. and Bravo, R. (Eds.), *Two-dimensional Gel Electrophoresis of Proteins*, Academic Press, New York 1984, pp. 3-36.
- [9] Garrels, J. I., *J. Biol. Chem.* 1979, 254, 7961-7977.
- [10] Anderson, N. L., *Two-Dimensional Electrophoresis, Operation of the ISO-DALT System*, Large Scale Biology Press, Washington, DC 1988, 162.
- [11] Bradford, M., *Anal. Biochem.* 1976, 72, 248.
- [12] Trace, R. P., Wold, L. E., Currie, L. M. and Young, D. S., *Clin. Chem.* 1982, 28, 890-899.
- [13] Merril, C. R., Goldman, D., Sedman, S. A. and Elbert, H. M., *Science* 1981, 211, 1437-1438.
- [14] Morrissey, J. H., *Anal. Biochem.* 1981, 117, 307-310.
- [15] Gard, D. L., Bell, P. B., Lazzarides, E., *Proc. Natl. Acad. Sci. USA*, 1979, 76, 3894-3898.
- [16] Matsumura, F., Lin, J.-C., Yamashiko-Matsumura, S., Thomas, G. P. and Topp, W. C., *J. Biol. Chem.* 1983, 258, 13954-13960.
- [17] Paulin, D., Forest, N. and Perreau, J., *J. Mol. Biol.* 1980, 144, 95-101.
- [18] Blobel, G. A., Moll, R., Franke, W. W., Kuster, K. W. and Gould, V. E., *Am. J. Pathol.* 1985, 121, 235-247.
- [19] Ochs, D. C., McConkey, H. E. and Guard, N. L., *Exp. Cell Res.* 1981, 135, 355-362.
- [20] Bhattacharya, B., Gaddamaniga, L. P., Valverius, E. M., Satomori, D. S. and H. L. Cooner, *Cancer Res.* 1990, 50, 2105-2112.
- [21] Sommer, C. L., Walker-Jones, D., Heckford, S. E., Worland, P., Valverius, A. Clark, R., McCormick, F., Stampfer, M., Abularen, S. and Gelmann, E. P., *Cancer Res.* 1989, 49, 4258-4263.
- [22] Trask, D. K., Bond, V., Zaichanski, D. A., Yawson, P., Suh, T. and Sager, R., *Proc. Natl. Acad. Sci. USA* 1990, 87, 2319-2323.
- [23] Trask, D. K., Bond, V., Zaichanski, D. A., Yawson, P., Suh, T. and Sager, R., *Proc. Natl. Acad. Sci. USA* 1990, 87, 2319-2323.





## LSB & LSP Information

Large Scale Biology Corporation

Large Scale Proteomics Corporation

---

### Large Scale Biology Corporation

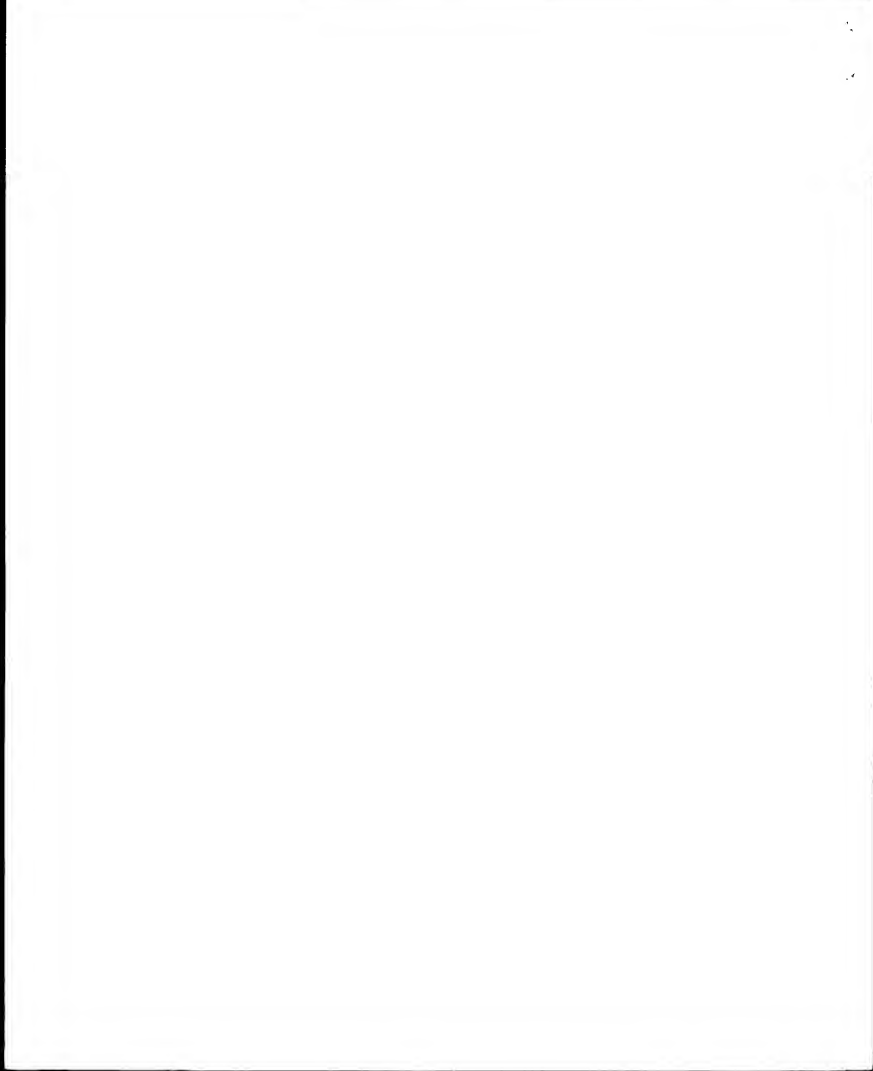
***Large Scale Biology Corporation is the leader in the integrated discovery, production and application of proteins - the functional units of all biological processes.***

Large Scale Biology Corporation (LSB, Vacaville, CA) and its subsidiary Large Scale Proteomics Corp. (LSP, Germantown, MD) are a biotechnology enterprise with the mission of accelerating the speed and productivity of the life sciences industry product discovery and development programs. Unique among biotechnology companies is LSB's integration of technologies to discover, analyze, manufacture and find new applications for proteins - the functional units of all biological processes.

Genomics companies have focused on deciphering genetic information, providing an initial but only partial understanding of biological processes. LSB's proprietary protein technologies can enable the transformation of genomic information into products such as drug targets, therapeutics, diagnostics for drug efficacy and toxicity, and traits for agricultural crops. Large Scale Biology has gone beyond the "genomics" realm in its business model and developed ways to integrate the discovery of gene function with quantitative protein analysis and protein manufacturing. This integration of technology platforms favorably positions LSB as a leading provider of valuable content to industry leaders in the fields of diagnostics, therapeutics, vaccines and agribusiness.

LSB was founded in 1987 with the goal of commercializing its proprietary GENEWARE viral vector system - a novel technology for gene expression. Using safe RNA viruses to transiently express genes in non-recombinant plants, LSB has positioned itself in the industry to provide cost-effective manufacturing and purification of diverse protein and peptide products. The same technology can be applied to the expression of libraries of foreign genes in an automated, high-throughput format to discover the function of genes with unparalleled efficiency. The GENEWARE system and associated proprietary technologies form the basis for LSB's functional genomics, biomanufacturing and a variety of proprietary products under development.

From its foundation, LSB understood the need to integrate functional genomic and protein manufacturing expertise with quantitative protein analysis and informatics to become a world-leader in the protein field. In 1999, LSB acquired a privately held pharmaceutical proteomics company originally founded in 1985. Large Scale Proteomics Corporation (a wholly





owned subsidiary of Large Scale Biology Corporation) is an industry leader in identifying and characterizing proteins in all types of biological samples for the discovery and development of new and more effective therapies, diagnostics, and agricultural products.

"Proteomics" is the study of the entire complement of proteins expressed in a cell, tissue, or organism. Proteomics can significantly improve drug discovery and development because most illness is associated with imbalances among, or malfunctions of, proteins. Only a small fraction of diseases can be attributed to the presence of a defective gene. Unlike classical genomics approaches that discover genes that may relate to a disease, LSP has developed a proprietary system called the ProGEx module for directly characterizing proteins associated with disease. Using this same technology, LSP can characterize the effects of candidate drugs intended to reverse a disease process, and to determine the degree to which this objective is achieved free of adverse side effects.

LSB and LSP have protected their many discoveries through an extensive portfolio of domestic and foreign patents and have developed commercial alliances and partnerships to exploit the value of their technologies. LSB and LSP scientists and engineers focus on the development and application of resources to help clients meet their objectives as well as the development of our own proprietary products for subsequent partnering with industry leaders.

A combined staff of 140 professionals operates from three locations in the United States, with a network of collaborators and affiliates throughout the US and Europe. Company headquarters, R&D laboratories and its Genomics division are located in Vacaville, California about 60 miles northeast of San Francisco. Process development and biomanufacturing take place in Owensboro, Kentucky, and LSB's Large Scale Proteomics Corporation subsidiary is located in Germantown, Maryland.

In August, 2000, LSB completed an initial public offering (IPO) of 5 million shares of common stock and now trades on the NASDAQ under the symbol LSBSC.

#### **Leadership - Large Scale Biology Corporation**

*Robert L. Erwin*, Chairman of the Board and Chief Executive Officer, founded LSB™ and has served as a director and officer since 1987. Mr. Erwin is the former chairman of the State of California Breast Cancer Research Council and currently serves on the University of California President's Engineering Advisory Council. He is Chairman of the Supervisory Board of Icon Genetics AG. As a co-founder of Sungene Technologies Corp., Mr. Erwin served as Vice President of Research and Product Development from 1981 through 1986. He has served on the Biotechnology Industry Advisory Board for Iowa State University. Mr. Erwin received his M.S. degree in Genetics from Louisiana State University and is an inventor on several LSB patents.

*David R. McGee, Ph.D.*, a co-founder of LSB and Senior Vice President and Chief Operating Officer, has been an officer since 1987. Prior to joining LSB, Dr. McGee was Vice President of Operations at Sungene Technologies Corporation from 1983 to 1987. Dr. McGee received his Ph.D. in Genetics from Louisiana State University and served as a faculty instructor of zoology and genetics at Louisiana State University.

*Laurence K. Grill, Ph.D.*, a co-founder of LSB and Senior Vice President, Research and Development, has served as an officer since 1987. Dr. Grill was the Manager of Plant Molecular Biology for Sandoz Crop Protection Corp. from 1984 to 1987 and Senior Research



Scientist in the Department of Molecular Biology at Zoecon Research Institute from 1980 to 1984. He received his Ph.D. from the University of California at Riverside with an emphasis on the molecular basis for viral gene expression in plants.

*R. Barry Holtz, Ph. D.*, Senior Vice President, Biopharmaceutical Manufacturing, has served the company as an officer since 1989 upon the acquisition of Holtz Bio-Engineering, which was founded in 1980. Dr. Holtz was a co-founder and Director of Research for MFI, Inc., the largest manufacturer of microencapsulated nutrients for agriculture and Director of Fundamental Research at Foremost-McKesson, Inc. Dr. Holtz received his Ph.D. in Biochemistry from Pennsylvania State University and served as Assistant Professor in the Department of Food Science and Nutrition at Ohio State University.

*Daniel Tusé, Ph.D.*, has been an officer of LSB since he joined the Company in 1995 as Vice President, Pharmaceutical Development. Dr. Tusé manages the company's pharmaceutical design and development programs, including LSB's novel vaccines and immunotherapeutics initiatives. Prior to joining LSB, Dr. Tusé was Assistant Director of SRI International's (Menlo Park, Calif.) Life Sciences Division. In his 17 years at SRI, Dr. Tusé developed extensive R&D experience in pharmaceuticals and specialty chemicals, serving an international list of clients. Dr. Tusé received his Ph.D. in Microbiology (1980, *cum laude*) with a minor in Toxicology from the University of California, Davis.

*John S. Rakitan*, a co-founder of LSB, Senior Vice President & General Counsel and Secretary, has served as an officer since 1988. Prior to joining LSB, Mr. Rakitan was an attorney in private practice. Mr. Rakitan received his J.D. degree from the University of Notre Dame.

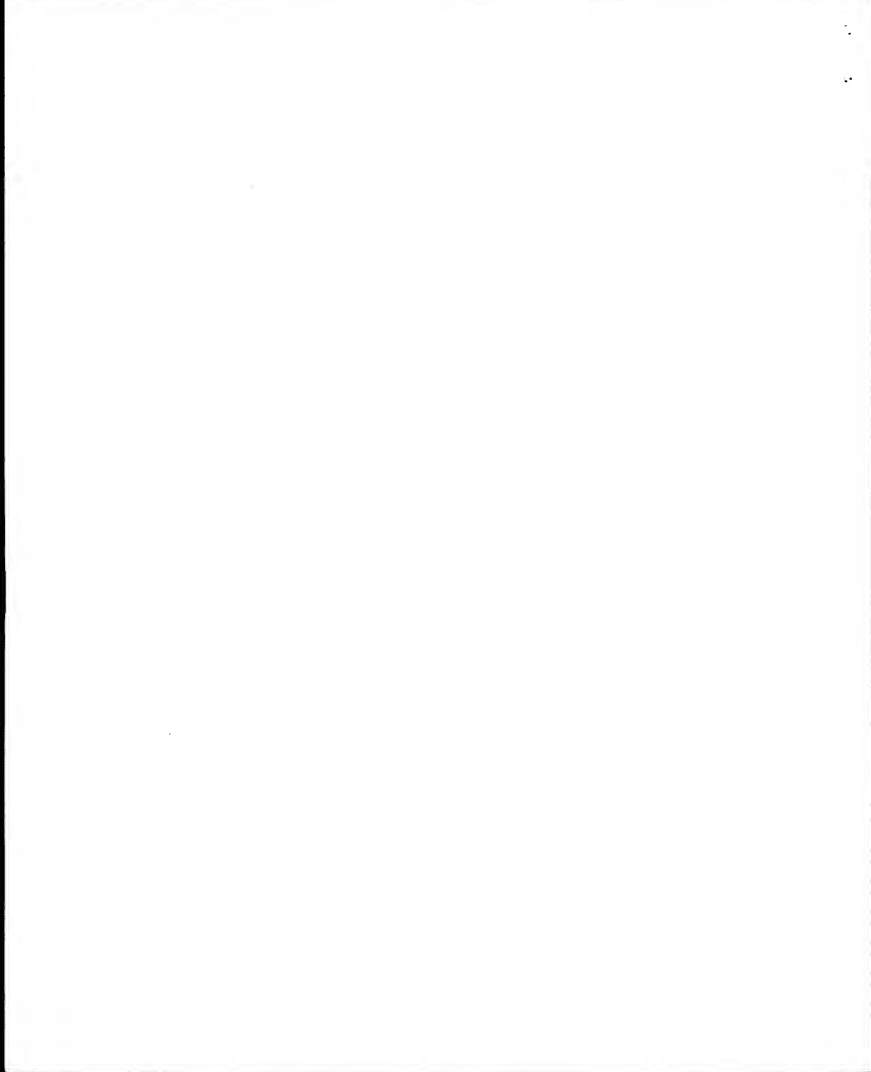
*Michael D. Centron*, Treasurer, has served as Controller since 1988 and was elected as Treasurer in 1991. Mr. Centron was Audit Supervisor for Varian Associates from June 1985 through July 1988, and he also worked for Arthur Young and Co. (currently Ernst & Young). Mr. Centron is a certified public accountant and received his M.B.A. degree from the University of California at Berkeley.

*Guy della-Cioppa, Ph.D.*, is an officer of the company and currently serves as Vice President, Genomics. Prior to joining the company in 1989, Dr. della-Cioppa worked for Monsanto Company in St. Louis, MO from 1984-1989 and was an NIH Postdoctoral Fellow at the Worcester Foundation for Experimental Biology in Shrewsbury, MA from 1983-1984. He received his Ph.D. in Biology from the University of California, Los Angeles.

*William M. Pfann* joined Large Scale Biology in August 2000 as Senior Vice President Finance and Chief Financial Officer. Mr. Pfann was formerly with PricewaterhouseCoopers LLP from 1969 to July 2000, most recently as the Risk Management Partner for the Western Region. He served in a number of management roles at PwC, including leader of the firm's Silicon Valley audit practice, National Director of the networking and communications sector and Managing Partner of the Northern California emerging business group, as well as Partner-in-Charge of the Oakland and Walnut Creek, California offices. Mr. Pfann received a B.S. degree from the University of California, Berkeley, in Business Administration and an MBA in Accounting from Golden Gate University.

#### [back to index](#)

© 2000 Large Scale Biology Corporation. All Rights Reserved Worldwide.



## Large Scale Proteomics Corporation

### Leadership - Large Scale Proteomics Corporation

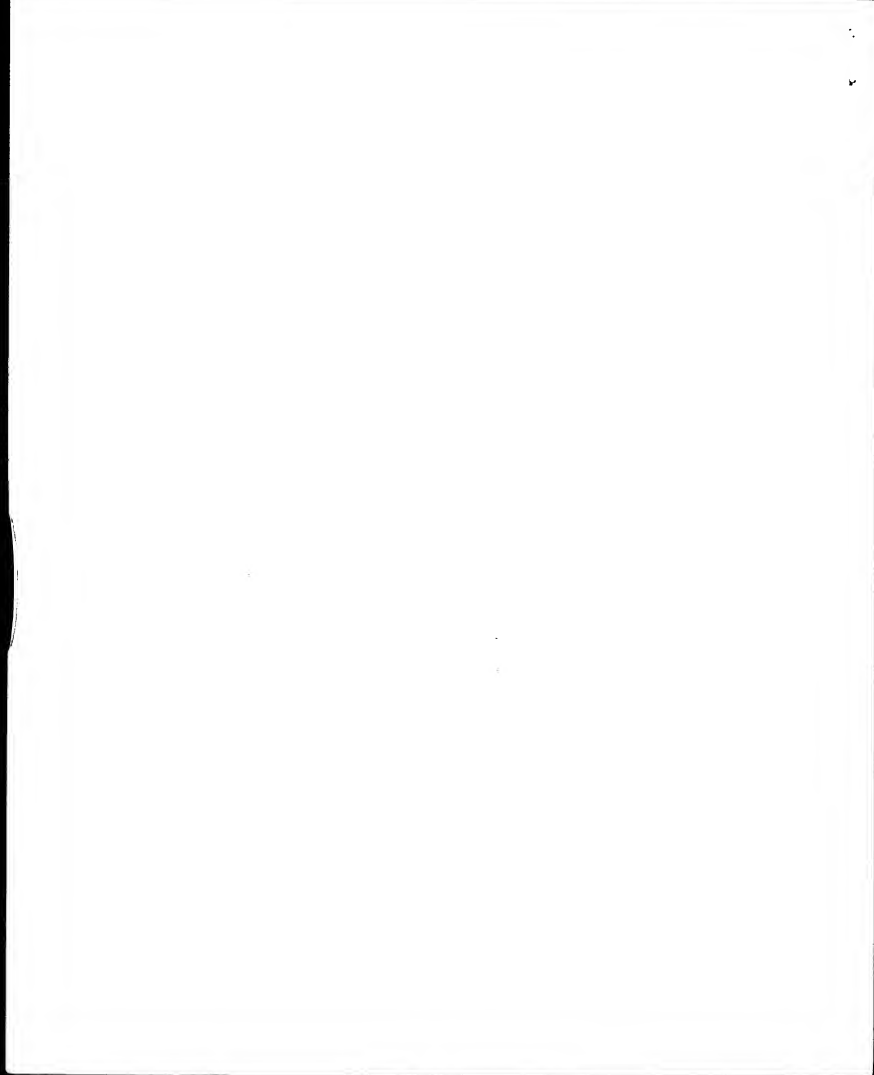
*N. Leigh Anderson, Ph.D.*, Chairman, President and CEO of Large Scale Proteomics Corporation (LSP™). Dr. Anderson obtained his B.A. in Physics with honors from Yale and a Ph.D. in Molecular Biology from Cambridge University (England) working with M. F. Perutz as a Churchill Fellow at the MRC Laboratory of Molecular Biology. Subsequently he co-founded the Molecular Anatomy Program at the Argonne National Laboratory (Chicago) where his work in the development of 2-dimensional electrophoresis (2-DE) and molecular database technology earned him, among other distinctions, the American Association for Clinical Chemistry's Young Investigator Award for 1982 and the 1983 Pittsburgh Analytical Chemistry Award. In 1985 Dr. Anderson co-founded LSP (originally Large Scale Biology Corp., Germantown, MD) in order to pursue commercial development and large-scale applications of 2-D electrophoretic protein mapping technology.

*Norman G. Anderson, Ph.D.*, Chief Scientist at LSP. Dr. Anderson has a distinguished record as an inventor. His career includes senior positions at Oak Ridge and Argonne National Laboratories (ORNL and ANL), more than 300 scientific publications, and the receipt of more than 20 prestigious awards in recognition of his work in science and technology. For his invention of the zonal ultracentrifuge, he received the John Scott Medal Award, and for the centrifugal fast analyzer, the Preis Biochemische Analytik für Klinische Chemie from Die Deutsche Gesellschaft für Klinische Chemie for the most outstanding analytical development in clinical chemistry worldwide during a 2-year period. In 1984 ANL awarded him its career patent leader award for the largest number of patents issued to an employee. At that time the commercial value of his inventions in terms of U.S. sales and royalties from foreign licensing were \$250 million and \$1 million, respectively. Dr. Anderson received his degrees at Duke University: a B.A. in Zoology, M.A. in Physiology, and Ph.D. in Cell Physiology. He holds 28 patents.

*Constance Seniff*, Vice President, Operations. Ms. Seniff has managed LSP's operations since 1993. Her background includes thirteen years in international business prior to joining LSP, five abroad in the employ of foreign firms. Ms. Seniff is responsible for helping formulate and implement business development and database commercialization strategies for LSP in coordination with the management of LSP's parent company, Large Scale Biology Corporation. Ms. Seniff has a B.Sc. degree in Business (with honors) from Florida State University.

*Robert J. Walden*, Vice President, Finance at LSP. Mr. Walden joined LSP in 1997 and has served as a director since 1999. He previously served as Vice President of Finance and Administration at Osiris Therapeutics, Inc., and as Chief Financial Officer at the American Type Culture Collection (ATCC). Mr. Walden received his degree in Finance from the University of Maryland.

*Jean-Paul Hofmann, Ph.D.*, Vice President, Software Development at LSP. Dr. Hofmann is a plant geneticist by training, having earned a B.S. in Biology, M.S. in Biochemistry and Genetics, and Ph.D. in Plant Genetics from the University of Orsay, Paris. He has extensive



experience in using 2-DE in agronomic research and in designing analytical software for 1- and 2-D applications. He has held senior scientific positions in industry and research institutes, in the U.S., France and the Ivory Coast.

*John Taylor, Ph.D.*, Vice President, Software Development and Bioinformatics. Dr. Taylor is the principal developer of Kepler™, LSP's analytical software for automated 2-DE pattern analysis. Prior to joining LSB, Dr. Taylor served as computer scientist in the Molecular Anatomy Program at Argonne, and on the research staffs of the University of Chicago and the Armed Forces Institute of Pathology in Washington, D.C. Dr. Taylor received a B.S. in Physics from the University of South Carolina, and a Ph.D. in Nuclear Physics from Duke University.

*Sandra Steiner, Ph.D.*, currently serves as Vice President Proteomics Applications. Prior to joining the Company, Dr. Steiner founded and directed the Molecular Toxicology Group at Novartis in Basel, Switzerland and was a member in several multi-disciplinary drug development project teams. Dr. Steiner received her Ph.D. in Toxicology/Pharmacology from the University of Basel, Switzerland.

**[back to index](#)**

© 2000 Large Scale Biology Corporation. All Rights Reserved Worldwide.

